# Name: Max Chamberlin, Candidate Number: MEC68

## Question 1: Hybrid Bisection Interpolation Algorithm

The bracket program is implemented in simple_bracket(), listed in the appendix:

Given an initial input, we search for a bracket along the direction of descent from the given input. We can think of the direction of descent as being like our positive x-axis to get our bearings. The input point then has a **negative** gradient along this descent direction. My algorithm searches for another point with a **positive** gradient along the direction of decrease. Once such a point has been found, we must have a bracket for the minimum. This is a consequence of the simple mathematical fact that for a continuous function, there must be a point x st: f'(x)=0 that is a minimum between a 'left' point with a negative gradient and the 'right' point with a positive gradient[1]

The search program is implemented in search(), listed in the appendix.

I used two different methods of interpolation in the search phase. The first method I implemented was a simple bisection. We update the bracket [L,R] for the minimum by interpolating midway between the bounds of the bracket. Then we select either [L,I] or [I,R] to form the new bracket. A minimum must be contained between a bracket where the leftmost point has a negative gradient and the rightmost point has a positive gradient, as previously explained; we select [L,I] or [I,R] based on this criterion.

However, I later adapted this method to produce a hybrid algorithm, which would give better convergence. If the function we hope to minimise is roughly parabolic near the minimum, then moving to the abscissa of a parabolic approximation should take us much closer to the true minimum - since smooth functions are well approximated by parabolas close to their minima. The parabola was fitting using gradient information. This typically resulted in faster convergence, but had the disadvantage that occasionally degenerate parabolas were created (ones were the interpolation point was at the same location as one of the previous bracket bounds). To remedy this problem, I opted to use a hybrid method that averaged between the midpoint and the parabolic interpolation point.

## Question 2: Testing the Algorithm

Below I plot the performance of the two algorithms described. I decided to initialise the points at various places of interest: for instance at a minima 0, for very large values of x, where the function oscillates much less and also for very small values of x. In the graphs, the bracket and initialisation point are given in the title. What can be seen is that both algorithms converge to the minimum. The hybrid algorithm I developed is faster than the midpoint algorithm, as can be seen from the relative steepness of the curves in the diagram below. Kinks in the curves are due to the fact that there are often several minima in a given bracket, and so the gradient of the interpolation point may not decrease as smoothly as we would expect. As bisection reduces the interval length by half on each iteration, the bisection method has precisely linear convergence. It remains to determine whether the hybrid algorithm has linear or super-linear convergence.

---

[1]**NB**: if the point is at an optima, my algorithm displaces the point slightly from the optima so that a direction of descent can be found, and a minimum can subsequently be bracketed.
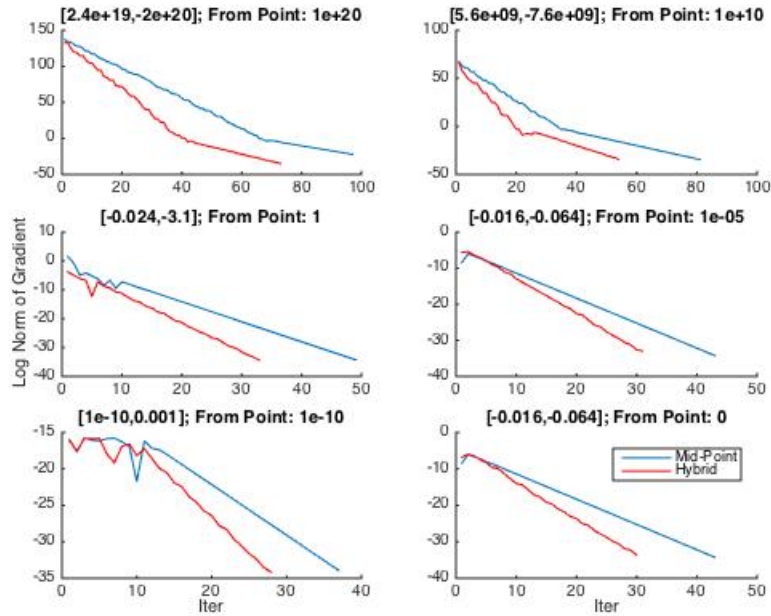
Figure 1: Convergence for Mid-Point and Hybrid Algorithm: Log Norm Gradient

I produced the plot below, which charts the distances between successive interpolation points. Since, the bisection method reduces the interval length by exactly half on each iteration, we should expect the completely linear rate of decrease as shown . The hybrid methods in some case looks like it may have reached a super-linear rate of convergence: for instance in the the graph from point 1e-10. However, most of the curves show seem to display a typical linear rate of decrease.
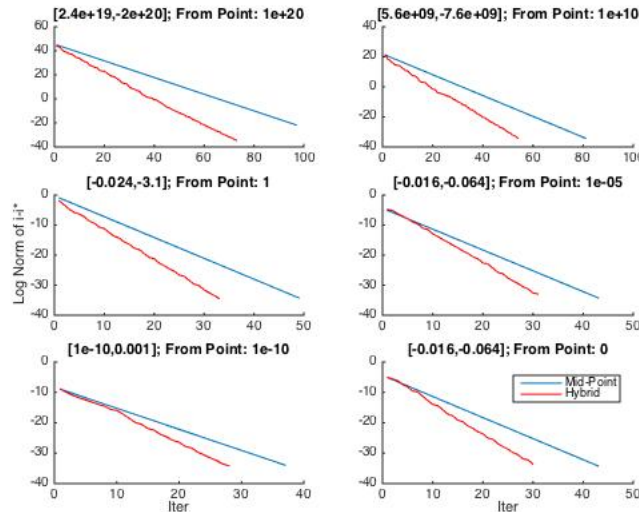


Figure 2: Convergence for Mid-Point and Hybrid Algorithm: Distances between Interpolation Points

Another interesting point to note, but that I decided not to produce graphs for, is that: since this function is completely symmetric about the 0-point, (and thus gradients and directions of descent are symmetric) my algorithm would produce exactly the same performance curves for negative initialisation points as for the corresponding positive points.

## Question 3

Since my line-search algorithm in question 1 only used gradient information to bracket and search for the minimum, a great many structural changes had to take place for the function to use the Wolfe conditions as a criterion for stopping. Ultimately, I implemented

the strong Wolfe conditions from the outline as described in Nocedal in Numerical Optimisation chapter 3, under: 'A Line Search Algorithm For the Wolfe Conditions' .This is given in the pseudo-code below. Zoom searches for the minimum, and LineSearch does the bracketing.

**Algorithm 3.5** (Line Search Algorithm).
Set $\alpha_0 \leftarrow 0$, choose $\alpha_{max} > 0$ and $\alpha_1 \in (0, \alpha_{max})$;
$i \leftarrow 1$;
**repeat**
    Evaluate $\phi(\alpha_i)$;
    **if** $\phi(\alpha_i) > \phi(0) + c_1\alpha_i\phi'(0)$ or $[\phi(\alpha_i) \geq \phi(\alpha_{i-1})$ and $i > 1]$
        $\alpha_* \leftarrow$ **zoom**$(\alpha_{i-1}, \alpha_i)$ and **stop**;
    Evaluate $\phi'(\alpha_i)$;
    **if** $|\phi'(\alpha_i)| \leq -c_2\phi'(0)$
        set $\alpha_* \leftarrow \alpha_i$ and **stop**;
    **if** $\phi'(\alpha_i) \geq 0$
        set $\alpha_* \leftarrow$ **zoom**$(\alpha_i, \alpha_{i-1})$ and **stop**;
    Choose $\alpha_{i+1} \in (\alpha_i, \alpha_{max})$;
    $i \leftarrow i + 1$;
**end (repeat)**

**Algorithm 3.6** (zoom).
**repeat**
    Interpolate (using quadratic, cubic, or bisection) to find
    a trial step length $\alpha_j$ between $\alpha_{lo}$ and $\alpha_{hi}$;
    Evaluate $\phi(\alpha_j)$;
    **if** $\phi(\alpha_j) > \phi(0) + c_1\alpha_j\phi'(0)$ or $\phi(\alpha_j) \geq \phi(\alpha_{lo})$
        $\alpha_{hi} \leftarrow \alpha_j$;
    **else**
        Evaluate $\phi'(\alpha_j)$;
        **if** $|\phi'(\alpha_j)| \leq -c_2\phi'(0)$
            Set $\alpha_* \leftarrow \alpha_j$ and **stop**;
        **if** $\phi'(\alpha_j)(\alpha_{hi} - \alpha_{lo}) \geq 0$
            $\alpha_{hi} \leftarrow \alpha_{lo}$;
        $\alpha_{lo} \leftarrow \alpha_j$;
**end (repeat)**

Figure 3: Strong Wolfe Line Search: Nocedal, Chapter 3

As a result of structural changes to the program, I decided to change the condition for bracketing from using derivative information to conditions on step lengths. According to Nocedal, we should bracket to ensure the step lengths bracket: $(\alpha_{i-1}, \alpha_i)$ contains step lengths satisfying the strong Wolfe conditions. This occurs if one of the following three conditions is satisfied: (i) $\alpha_i$ violates the sufficient decrease condition; (ii) $(\alpha_i)$ $(\alpha_{i-1})$; (iii) $(\alpha_i)$ $0$ - where $\phi(\alpha) = f(x + dir \cdot \alpha)$; $\phi'(\alpha) = f'(x + dir \cdot \alpha)$. These changes are also reflected in my code, for when one of those conditions is found to hold we undertake the search that leads us to a minimum.

To implement the algorithm above neatly, I decided that all the information relating to a step-size should be held in one place. As such, I created a 'step size' class that stored all the information related to different step sizes. This included: the locations stepped to; the vector x they were stepping away from; the descent directions, their gradient and function values at the current location; and their gradients along the descent direction. I also gave each step object a counter to count their function evaluations.

Interpolation: because the algorithm was previously constructed on the basis of derivative information, which is no longer possible in this setting, it would not be possible to implement the simple parabolic interpolation as presented in question 1. The bracket no longer contains a negative and positive gradient for parabolic interpolation. Consequently, I decided to use the simple bisection method that was initially discussed and **tested** for in question 2. However, I have also included a modified hybrid bisection algorithm which interpolates based on different principles from those expounded in question 1.

## Question 4:

The code to perform steepest descent and conjugate gradients is listed within the appendix. My programs perform an inexact line search using the line search method that had been developed in question 3. In this question we use bisection to interpolate.

In the diagrams below, I plot the trajectories of the respective functions for the point initialised at (x,y)=[-1,0.5].
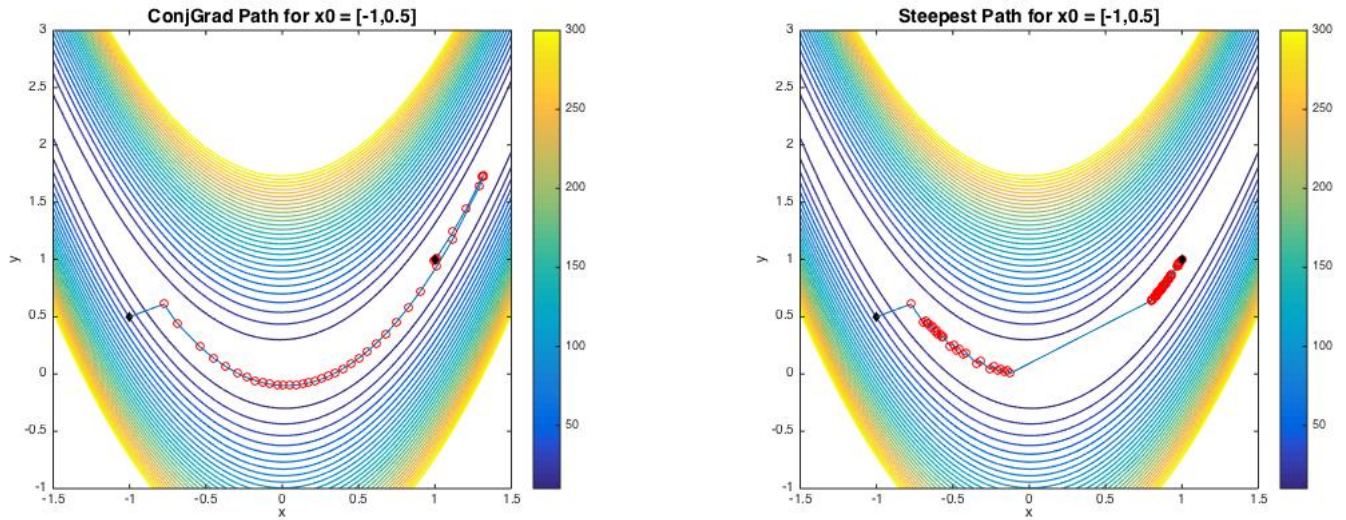
Figure 4: Paths over Rosenbrock Valley for SD and CG Algorithms

From the of the respective paths above, we can see that the conjugate gradients takes a smoother route along the valley floor, following a path that makes larger jumps between successive sampled points. The steepest descent method tends to spend a lot of time moving perpendicular to the contour lines of the figure, making very little progress towards the true minimum as it follows the paths of steepest descent. Because the directions of conjugate gradients are somewhat orthogonal, there is much less wasted movement as it progresses towards the global minimum.

Below, I have also tabulated the mean function calls and main-loop iterations taken for the respective algorithms to converge to the minimum, with initial points sampled from 0 to 2 randomly in both x and y. I averaged over 30 values. The criterion for convergence is that the norm of the gradient must be less that 1e-6. As we can see, on average Steepest Descent takes 1000 times more function calls and 100 times as many main loop iterations as conjugate gradients to converge to the minimum. But why should SD perform worse when scrutinised by function call than by iteration? This is probably due to the fact that it's very difficult to find a small enough step-size with a sufficient decrease along a descent direction when we approach the minimum, since descent directions are almost parallel to one another. This makes progress increasingly difficult. In contrast, CG, which looks for orthogonal directions, is not blighted by the same problem and can easily find large step sizes between successive points .

| Tolerance =1e-6 | Conjugate Gradients | Gradient Descent |
|---|---|---|
| Mean f Calls | 970.3000 | 129,230 |
| Std Deviation f Calls | 910.8231 | 50,989 |
| Mean Iters | 42.4000 | 4,562.1 |
| Std Iters | 36.5547 | 1,819.4 |

I have also graphed the results of running the algorithm on a test point, so that convergence can be seen visually. As can be seen, conjugate gradients converges much more rapidly than steepest descent. Both algorithms tend to converge to the minimum when initialised fairly close to the global minimum. However, whilst Conjugate gradients seems to converge even from afar (perhaps due to the robust search directions), the steepest descent algorithm seems to cycle aimlessly through a set of points, taking smaller and smaller steps towards a minimum.
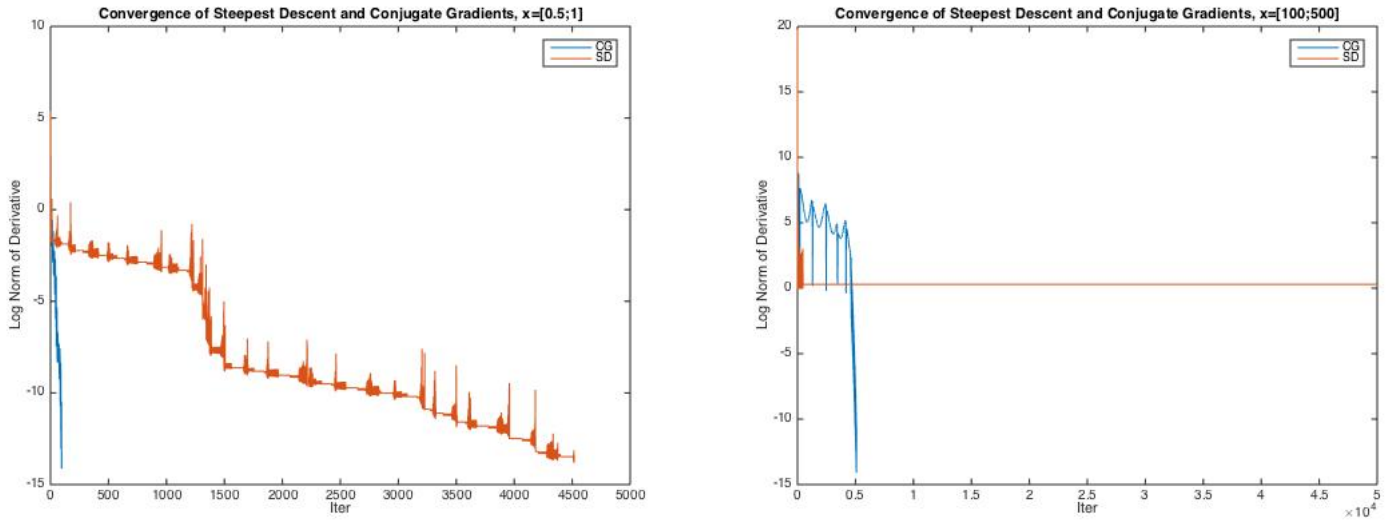
Figure 5: Graph of Convergence for Conjugate Gradients and Steepest Descent

**Note:**

The Rosenbrock function can be trivially solved by applying a change of variables: $z_1 = (1-x)$, $z_2 = 10*(y-x^2)$. In which case $f(x,y) = g(z_1, z_2) = \mathbf{z^T I z}$. This latter function can be solved with an exact line search, which would be much less computationally expensive. Because the shape of the quadratic cost is uniform, steepest descent would solve the problem in 2 iterations. However, I take it that the question was asking for a robust method to measure the performance of the algorithms, rather than exploiting other mathematical manipulations.

## Question 5:

The programs to perform an exact line-search Steepest Descent and Conjugate Gradients are listed in the appendix.

In general the B matrices are better suited to gradient descent than the A matrices: the contours of their objective function are more 'spherical' than they are like 'ellipsoids'. This can be seen from the tabulation of their eigen-values below. The largest eigen-value in B is only slightly greater ( by a factor of 10) than the smallest eigen -value. In the A matrices, the largest eigen-value is more than 10,000 times as great as the smallest. This will indeed produce a very narrow elliptical objective function. As a result, we should expect that steepest descent would perform much more poorly than conjugate gradients on the B matrices, as it gets stuck in the narrow valleys (see figure below).

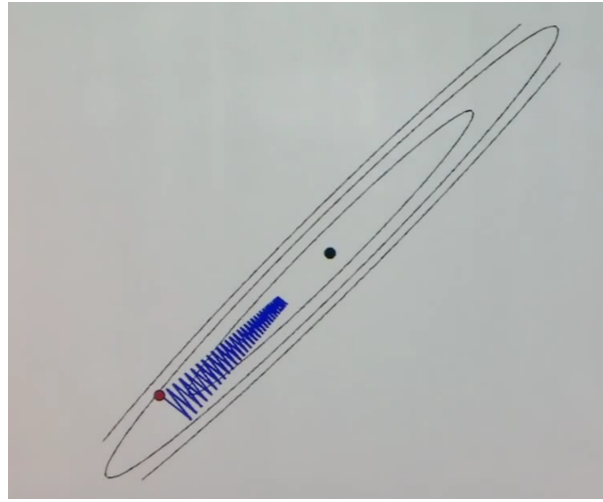| matrix/eigen_values | $\lambda_1$ | $\lambda_2$ | $\lambda_3$ | $\lambda_4$ | $\lambda_5$ | $\lambda_6$ | $\lambda_7$ | $\lambda_8$ | $\lambda_9$ | $\lambda_{10}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| matrix A10 | 1.0825 | 1.0460 | 0.1082 | 0.1046 | 0.0105 | 0.0108 | 0.0011 | 0.0010 | 0.0001 | 0.0001 |
| matrix B10 | 1.0266 | 1.0267 | 10.4256 | 10.4256 | 10.2667 | 10.2667 | 1.0267 | 1.0426 | 1.0426 | 1.0267 |

Figure 6: An example of Steepest Descent stuck in a Narrow Ellipse, From: tumblr_naubqg3kYK1qc38e9o1_1280.png

Below I have plotted graphs showing the performance of both SD and CG from or an initial point x, whose elements are identically set to 1, and with the b vector in the $b^T x$ term of the objective function randomised as suggested. As can be seen, for the more elliptical A matrices, CG tends to reach the minimum thousands of times faster. SD does converge, however the rate is rather slow- with the number of iterations until convergence thousands of times larger than the dimensions of A10. For the more spherical matrices, CG still outperforms steepest descent, however the performance of the two methods is much closer.
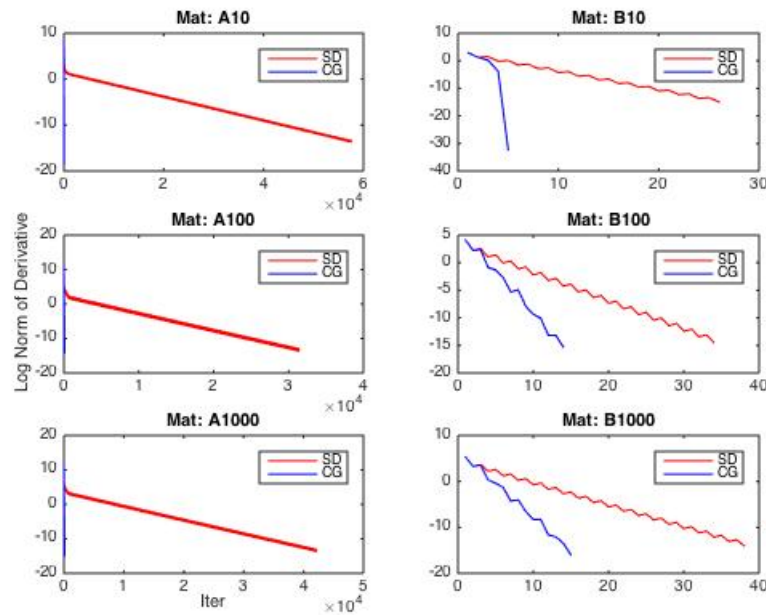


Figure 7: Convergence of the Conjugate Gradients and Steepest Descent over A,B matrices

To produce a more general set of results, I also tabulated an average number of function calls until convergence- specified by the tolerance 1e-4. Another interesting point to note is that the standard deviation of the number of function calls for conjugate gradients is very low, often 0. From a theoretical perspective, conjugate gradients should find the exact minimum after a number of evaluations roughly equal to A's or B's dimension (Similar to how it takes less than N steps to find the minimum of an N dimensional matrix by moving along the direction of it's eigenvectors). However, to reach the tolerance threshold, it takes a little more time for 2 A matrices and much less time for the B matrices. From the table below, we can see that SD performs only marginally worse (by a factor of 2 or 3 times the number of f calls) for fairly spherical matrices, but many thousands of times worse for elliptical matrices.

| Tolerance =1e-4 | A10 | B10 | A100 | B100 | A1000 | B1000 |
|---|---|---|---|---|---|---|
| CJG Mean f Calls | 15.8 | 5 | 94.53 | 14 | 136 | 15 |
| GJG Std Deviation f Calls | 0.43 | 0 | 2.74 | 0 | 0 | 0 |
| SD Mean f Calls | 10,930 | 22.6 | 15,809 | 34 | 17,018 | 38 |
| SD Std Deviation f Calls | 1,343 | 3.32 | 272.74 | 0 | 91.28 | 0 |

## Question 6:

For this question, I used the same algorithms used in question 4, but supplied with arguments for the objective and gradient functions of the A and B matrices. The graph below shows the performance of the two algorithms from or an initial point x, whose elements are identically set to 1 and with b in the $b^T x$ term randomised.
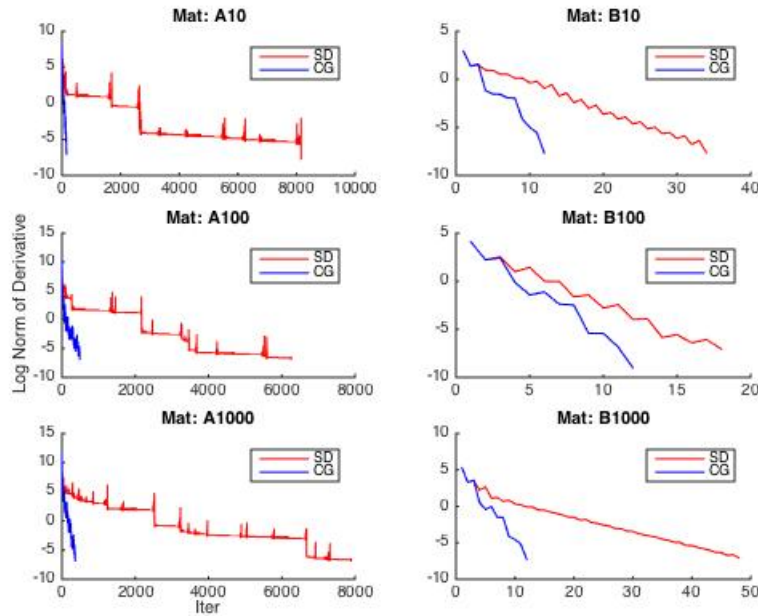


Figure 8: Convergence of the Conjugate Gradients and Steepest Descent over A,B matrices

The graphs show that rate of convergence is generally much less smooth for the A matrices, which is perhaps to be expected since the method uses an inexact line search. This is because when the contours of the objective function are plotted it can be shown to have a very narrow elliptical shape, and as such even small movements along a direction of descent can result in much larger changes in gradient. The graphs of convergence for the B matrices are much smoother because the ellipses are not generally moving into a 'valley' of the 'ellipse', but more in away which is parallel to the contours. This results in a much smaller change in gradient.

I tabulated the average number of iterations and function calls until convergence below. Generally, conjugate gradients performed comparably as well as before for the more spherical objective functions as defined by the B matrices. However, Conjugate gradients took much longer (10 times as long) to converge with the more elliptical A matrices. With an exact line search, we had a theoretical guarantee that convergence would take place in less than n iterations (where n is the dim of A or B), with an inexact line search there is no such guarantee and convergence takes considerably longer.

With Steepest Descent, there seemed to be a reduction in the number of calls to the line search function, by a factor of around 2/3. This is to be expected since an inexact line search should takes us further away from a narrow valley within an ellipse, than an exact line search. This means that we can move more closely to the contours of the ellipse as conjugate gradients does, rather than getting stuck iterating between the contours of a narrow ellipse. However, if we look at the number of function calls in their totality (rather than comparing calls to line search within the main loop) the number of function evaluations has actually increased. This can be seen in the second table below.

Main Loop Iters:

| Tolerance =1e-4 | A10 | B10 | A100 | B100 | A1000 | B1000 |
|---|---|---|---|---|---|---|
| CJG Mean iters | 187.7333 | 16.10 | 258.2667 | 17.0667 | 288 | 19.6667 |
| GJG Std Deviation iters | 83.6095 | 3.0212 | 26.8057 | 1.7991 | 87.7765 | 1.2411 |
| SD Mean f iters | 3,376 | 17.2000 | 2174 | 14.5000 | 2528 | 39.6000 |
| SD Std Deviation f iters | 327 | 6.8097 | 198 | 4.9532 | 237 | 9.6369 |

Function Calls:

| Tolerance =1e-3 | A10 | B10 | A100 | B100 | A1000 | B1000 |
|---|---|---|---|---|---|---|
| CJG Mean f Calls | 3,833 | 146.9 | 5.2236e+03 | 162.4333 | 5.5659e+03 | 185.8667 |
| GJG Std Deviation f Calls | 1,993 | 54.73 | 540.0611 | 29.1366 | 2.0026e+03 | 18.5300 |
| SD Mean f Calls | 17,330 | 214.8667 | 39,659 | 163.4333 | 43,975 | 591.6000 |
| SD Std Deviation f Calls | 4,565 | 93.9269 | 3,872 | 62.7764 | 5,667 | 150.8817 |

# Appendix:

**Q1: simple_bracket()**

```matlab
function [left, right, g_eval]=simple_bracket(gradF, x)
% INPUT
% gradF  - gradient of function
% x      - initial point
% OUTPUT
% left   - left bracket
% right  - right bracket
% g_eval - gradient evaluations


% Preliminaries: Find a descent direction
%-------------------------------------------------------
g   = gradF(x);
dir = -sign(g);
% Code below handles rare 'bad inputs', eg saddle points.
% i1 counts the number of function calls
i1=1;
while(dir==0 && i1<10)
    % The displacement increases geometrically to escape saddle points.
    x=x+(1e10^i1)*eps;
    g   =    gradF(x);
    dir=-sign(g);
    i1=i1+1;
end
% Did we escape, do we have a descent direction?
if (dir==0)
    error('Could not escape optima');
end
%-------------------------------------------------------




% Main Algorithm:
%-------------------------------------------------------
alpha              = 0.1;       % Initial step length
fac                = 4;         % Increase factor
alpha_prev         =0;          % Prior step length
```

```
39 i2                        =0;            % Counts gradient calls of main loop
40
41            % Main Loop:
42            %————————————————————————————————————
43            % Terminate on 1000th function call
44            while(i2<1000)
45                 % Take a step in the descent direction
46                 g_alpha=gradF(x+alpha*dir);
47                 i2=i2+1; g_eval=i2+i1;
48                 % Do we have a bracket for the minimum?
49                 if g_alpha*dir > 0
50                      right=x+alpha*dir;
51                      left=x+alpha_prev*dir;
52                      return;
53                 end;
54                 %Increase step sizes
55                 alpha_prev=alpha;
56                 alpha = fac*alpha;
57            end
58            error('error: Search exceeds 1000 iterations');
59 %————————————————————————————————————————————
```

## Q1: search algorithm()

- Both bisection and hybrid methods performed

```
1  function [a,interp,b,numEvals,gradDiff, Diff] = search(gradF,a,b,type)
2  %[a,x,b,h,i]=simple_bracket(@f,@gradF,1e19,@bisect)
3  % INPUT
4  % gradF - gradient of function to be minimised
5  % a      - left bracket
6  % b      - right bracket
7  % type   - 'hybrid' for hybrid interpolation, otherwise mid point used
8
9  % OUTPUT
10 % a        - updated left bracket
11 % b        - updated right bracket
12 % interp   - interpolation point* estimate for minimum
13 % numEvals - number of gradient evaluations
14 % gradDiff - log norm of successive gradients - for plotting
15 % Diff     - absolute difference between interps - for plotting
16
17
18 %Preliminaries
19 %---------------------------------------------------------------------
20 % tolerance for minimum acceptance
21 tol=eps*1e1;
22
23 % gradient of end points
24 gradA=gradF(a);
25 gradB=gradF(b);
26 % interpolation point
27 interp=(a+b)/2;
28 % number of function evaluations so far
29 numEvals=2;
30 % Store absolute difference in gradient to plot
31 gradDiff=nan(500); gradI=nan;
32 Diff=nan(500);
33
34
35 % Main Loop
36 %---------------------------------------------------------------------
37 % Bisect until convergence  within tolerence or
38 while (abs(a-b)>tol && numEvals <1000)
39     prev=interp;
40     %1. Interpolate between the bracket
41     %-------------------------------------------------
42     % take mid point as default interpolation
43     interp=(a+b)/2;
44
45     % update using hybrid or mid point criterion
46     if strcmp(type,'hybrid')
47         % construct a parabola y=px^2+qx+d through a and b, from deriv info
48         p=0.5*(gradA-gradB)/(a-b);
49         q=0.5*(gradA+gradB-2*p*(a+b));
50         % estimation of min point of parabola
51         min=-0.5*q/p;
```

```
52          % hybrid estimation offsets averages parabolic with mid point
53          interp=((a+b)/2 +min)/2;
54      end
55
56
57      %2. Store the absolute difference in gradients to plot
58      %----------------------------------------------------
59       gradIPrev=gradI;
60       gradI=gradF(interp);
61       gradDiff(numEvals)=abs(gradIPrev-gradI);
62       Diff(numEvals)=abs(prev-interp);
63
64      %3. Update the bracket of the minimum
65      %----------------------------------------------------
66      if (gradI==0)
67          return
68      % Does the interpolation point have the same gradient as A
69      elseif (sign(gradA)==sign(gradI))
70          a=interp; gradA=gradI;
71      % If not it must have the same sign as B
72      else
73          b=interp; gradB=gradI;
74      end
75      numEvals=numEvals+1;
```

## Question 3: Line Search

### Step object

```
1  classdef step  < handle
2      properties
3          size; %step size
4          x;   % step from where
5          f =[]; % function eval
6          df =[];% derivative information
7          dir; % direction normed
8          count =0;
9          loc; %location after step
10         dfdir=[]; % gradient in line direction
11     end
12
13     methods
14         % initialise the step size values
15         function init(obj,size,x,dir)
16             obj.size=size;
17             obj.x=x;
18             obj.dir=dir/norm(dir);
19             obj.loc=obj.x+obj.dir*obj.size;
20             obj.f=[];
21             obj.df=[];
22         end
23         % evaluate the function at the step's location
24          function evalF(obj,f)
25           if (isempty(obj.f))
26             obj.f=f(obj.loc);
27             obj.count=obj.count+1;
28           end
29          end
30         % initialise the gradient
31          function evalDF(obj,df)
32           if (isempty(obj.df))
33             obj.df=df(obj.loc);
34             obj.count=obj.count+1;
35             obj.dfdir=dot(obj.df,obj.dir);
36           end
37          end
38          % copy values from another step, but keep function evaluation counter the same
39           function obj = copy(obj,obj2)
40             obj.size=obj2.size;
41             obj.x=obj2.x;
42             obj.dir=obj2.dir;
43             obj.loc=obj2.loc;
44             obj.df=obj2.df;
45             obj.f=obj2.f;
46             obj.dfdir=obj2.dfdir;
47           end
48     end
49
50  end
```

### Line Search() with Zoom() nested

```matlab
function [opt,z] = LineSearch(f,gradF,x,dir, type)

% implementation of wolfe line search
%
% f        - function handle of objective function
% gradF    - function handle of gradient
% x        - current iterate
% dir      - search direction must be a descent
% a        - the a vars are step objects
% type     - type of interpolation

% Output
% opt      - min point
% z        - number of function evaluations


a0 =step; ai=step; a_pre=step;

% c1= suff dec, c2= curvature condition NB: c1 <= c2
%c1 = 0.01; c2 = 0.1;
c1 = 0.01; c2 = 0.1;

a0.init(0,x,dir); a0.evalF(f); a0.evalDF(gradF);
ai.init(1,x,dir); a_pre=a_pre.copy(a0);

iter = 1;
while (iter<1000) %cap on iterations is 1000
    %'main'
    ai.evalF(f);
    % check if current iterate violates sufficient decrease
    if  ai.f > a0.f + c1*ai.size*a0.dfdir || (ai.f >= a_pre.f && iter > 1)
        %  Acceptable point between a_pre and a_i because (c1 > c2)
        [opt,z] = nocZoom(a_pre, ai);
        z=z+a0.count+ai.count+a_pre.count;
        return;
    end

    ai.evalDF(gradF);
    % current iterate has sufficient decrease, but are we too close?
    %if(ai.dfdir >=c2*a0.dfdir)
     if abs(ai.dfdir) <=-(c2*a0.dfdir)
        % Wolfe fullfilled, quit
        opt= ai.loc;
        z=a0.count+ai.count+a_pre.count;
        return;
    end
    % are we ahead of the minimum?
    if (ai.dfdir >= 0)
        % there has to be an acceptable point between a_pre and ai
        [opt,z] = nocZoom(ai,a_pre);
        z=z+a0.count+ai.count+a_pre.count;
        return;
    end
    a_pre=a_pre.copy(ai);
    %Double ?i step size ;
    ai.init(ai.size*2,x,dir);
    iter = iter + 1;
end



function [m,z] = nocZoom(aLo,aHi)
% Interpolated step size

aj=step;
aHi.evalDF(gradF);
aLo.evalDF(gradF);

  while 1
      'zoom';

  %update mid point with parabolic interpolation if different criterion used
    if strcmp(type,'hybrid')
        % construct parabola
        aLo.evalF(f);aHi.evalF(f);
        num = aLo.dfdir(aHi.size-aLo.size);
        denom = 2*(aLo.f+aLo.dfdir*(aHi.size-aLo.size)-aHi.f);
        % estimation of min point of parabola
        min= aLo.size+num/denom;
        % hybrid estimation offsets averages parabolic with mid point
        size=((aHi.size+aLo.size)/2 +min)/2;
    else
        size=(aHi.size+aLo.size)/2;
    end

     aj.init(size,x,dir);
```

```
88      %aj.evalDF(gradF);
89      aj.evalF(f);
90
91      if aj.f > a0.f + c1*aj.size*a0.dfdir || aj.f >= aLo.f
92          % No sufficient decrease for a => a= maximum of the interval
93          aHi = aHi.copy(aj);
94      else
95          aj.evalDF(gradF);
96          % strong wolfe fullfilled?
97          %if aj.dfdir >=c2*a0.dfdir
98          if abs(aj.dfdir) <=-(c2*a0.dfdir)
99
100             m=aj.loc; z=aj.count;
101             return;
102         end
103         % if slope positive and aHi > aLo
104         if aj.dfdir*(aHi.size-aLo.size) >= 0
105             aHi = aHi.copy(aLo);
106             %aLo = aLo.copy(aj) ;
107         end
108           aLo=aLo.copy(aj);
109      end
110
111    end
112
113 end
114
115 end
```

## Question 4

**Steepest descent()**

```
1  function [x,logAbsG,evals,k]=steepest(x, gradF, f)
2  % Input:
3  % x       -  initial guess of the solution
4  % gradF -   gradient of function
5  % f       -  function to minimise
6
7  % Output:
8  % x            -  minimum of f
9  % logAbsG      -  the log of the norm of the gradient over iters
10 % evals        -  number of function evals
11 % k            - nuumber of iterations
12
13 %Output:
14
15 % tolerance for stopping algorithm
16 tol = 1e-4;
17 % iteration number
18
19 k = 1;
20 % maximum number of iterations
21 MAX=5000;
22
23 % g- gradient at x; x next location; p- conjugate gradient
24 g = gradF(x);
25
26 % num evals is 0
27 evals=0;
28 % log  norm of the gradient for plotting
29 logAbsG=nan(MAX,1);
30 logAbsG(k,1)=log(norm(g));
31
32 while norm(g) >tol & k<MAX
33     % get new point guess from minimum and num evals from Line Search
34     % search in direction of steeopest descent
35     [x,evals1] =  LineSearch(f,gradF,x,-g);
36     % update number of evaluations
37     evals=evals+evals1;
38     % get new gradient
39     g=gradF(x);
40     % update number of evaluations
41     evals=evals+1;
42     % update iter
43     k=k+1;
44     % store gradient info for plotting
45     logAbsG(k,1)=log(norm(g));
46 end
47 end
```

**Conjugate Gradients()**

```matlab
function [xk,logAbsG,evals,k]= conjGrad(x1, gradF, f)
% minimise along conjugate vectors: pk which are conjugate to A
% At each step we chose ak by an inexact line searc
% Input:
% x      -  initial guess of the solution
% gradF  -  gradient of function
% f      -  function to minimise

% Output:
% x          -  minimum of f
% logAbsG    -  the log of the norm of the gradient over iters
% evals      -  number of function evals
% k          -  nuumber of iterations

%Output:


% tolerance for stopping algorithm
tol = 1e-4;
% iteration number

k = 1;
% maximum number of iterations
MAX=5000;

% g- gradient at x; x next location; p- conjugate gradient
gk = gradF(x);

% search direction is -grad
pk=-gradF(x);

% num evals is 0
evals=0;

% log  norm of the gradient for plotting
logAbsG=nan(MAX,1);
logAbsG(k,1)=log(norm(g));

while norm(gk) >tol & k<max

            % Check conjugate gradient is a descent direction
            if gradF(xk)'*pk>0
                 pk=-pk;
            end;
            % Perform line search
            [xk_next,evals1] = LineSearch(f,gradF,xk,pk, 'hybrid');
            % Update number of ealuations
            evals=evals+evals1;

            % This section updates the conjugate gradients
            % according to the method used in Barber, BRML
            gk_next=gradF(xk_next);
            evals=evals+1;
            bk =(gk_next'*gk_next)/(gk'*gk);
            pk_next=-gk_next+bk*pk;
            k=k+1; pk=pk_next; gk=gk_next; xk=xk_next;

            % we also apply the heuristic that the gradients are reset
            % every 100 iterations
            if floor(k/100)==1
                 pk = -gk;
                 bk=0;
            end;

            % update logAbs G for plotting
            logAbsG(k,1)=log(norm(gk));
end
end
```

## Question 5 and 6

**Steepest descent()**

```matlab
function [x,k,logAbsG]=steepestMat(x,A,seed)
% Input:
% x      -  initial guess of the solution
% A      -  A objective function
% note b is created with a random seed

% Output:
% x          -  minimum of f
% logAbsG    -  the log of the norm of the gradient over iters
```

```matlab
10  % k              - nuumber of iterations
11
12
13  % tolerance for stopping algorithm
14  tol = 1e-6;
15
16  % Create Objective vector
17  rng(seed)
18  N=length(A);
19  b=2*rand(N,1)-1;
20
21  % g- gradient at x
22  g = A*x-b;
23  % negative of gradient
24  p = -g;
25
26  % number of iterations
27  k=1;
28
29  % Maximum number of iterations
30  MAX=10000;
31  % log Abs gradient for plotting
32  logAbsG=nan(MAX,1);
33  logAbsG(k,1)=log(norm(g));
34
35  % whle not converged do:
36  while norm(p) >tol & k<MAX
37
38          % exact step size
39          a= g'*g/(p'*A*p);
40          % update x
41          x=x+a*p;
42          % update the search directions and conj grad
43          g=A*x-b; p=-g;
44          k=k+1;
45          % store logabsG for plotting
46          logAbsG(k,1)=log(norm(g));
47  end
48
49  end
```

## Conjugate gradients

```matlab
1   function [x,k,logAbsG]=conjGradMat(x,A,seed)
2   % Input:
3   % x      -  initial guess of the solution
4   % A      -  A objective function
5   % note b is created with a random seed
6
7   % Output:
8   % x              -  minimum of f
9   % logAbsG        -  the log of the norm of the gradient over iters
10  % k              - nuumber of iterations
11
12
13  % tolerance for stopping algorithm
14  tol = 1e-6;
15
16  % Create Objective vector
17  rng(seed)
18  N=length(A);
19  b=2*rand(N,1)-1;
20
21  % g- gradient at x
22  g = A*x-b;
23  % conjugate gradient
24  p = -g;
25
26  % number of iterations
27  k=1;
28
29  % Maximum number of iterations
30  MAX=10000;
31  % log Abs gradient for plotting
32  logAbsG=nan(MAX,1);
33  logAbsG(k,1)=log(norm(g));
34
35  % whle not converged do:
36  while norm(p) >tol & k<MAX
37
38          % The method is given by Barber, BRML
39          % Exact Line Search
40          a= g'*g/(p'*A*p);
41          % Update x
42          x_next=x+a*p;
43          % Update grad, conj
44          g_next=A*x_next-b;
```

```matlab
            b =(g_next'*g_next)/(g'*g);
            p_next=-g_next+b*p;
            k=k+1; p=p_next; g=g_next; x=x_next;

            % atore log nor grad for plotting
            logAbsG(k,1)=log(norm(g));
    end

end
```