

# MATH 6644

## Homework 2

Stephan Boettcher

April 14, 2015

### Question 1

*Can the performance of the Newton iteration be improved by a linear change of variables? That is, for nonsingular  $N \times N$  matrices  $A$  and  $B$ , can the Newton iterates for  $F(x) = 0$  and  $AF(Bx) = 0$  show any performance difference when started at the same initial iterate? What about the chord method?*

A linear change of variables is a technique used to reduce a difficult to a simpler one. This is commonly done by substituting values or expressions for ones that depend on other variables. By reducing the number of dependent variables in a set of expressions, the corresponding  $N \times N$   $A$  matrix becomes more sparse. For a given Newton iterate,  $F(x) = 0$ , the nonsingular  $A$  and  $B$  matrices can be used to modify the sparsity of the iterates, given by  $AF(Bx) = 0$ , to make them easier to evaluate. This would be beneficial to Newton's method as a more sparse iterate makes evaluating the Jacobian less expensive. With Newton's method, both the iterate,  $AF(Bx_n)$  and the corresponding Jacobian are evaluated with each iteration. Thus, if they can be made less costly to evaluate, a performance increase can be expected. However, the degree of sparsity and resulting performance increase must offset the cost of the extra matrix-vector multiplications, which may quickly and completely offset the gains mentioned above. The other way the performance of Newton's method could be expected to increase is from a reduction in the number iterations necessary to converge. In order to evaluate this, we must first look at Newton's equation.

Newton's Method converges quadratically to  $x$  and the error of the function is governed by the equation:

$$\|\vec{x}_{n+1} - \vec{x}_n\| = \|(F'(\vec{x}_n))^{-1}F(\vec{x}_n)\|$$

Substituting in  $A$  and  $B$  into the above equation, we get:

$$\|\vec{x}_{n+1} - \vec{x}_n\| = \|(ABF'(B\vec{x}_n))^{-1}AF(B\vec{x}_n)\|$$

Since we know that the  $A$  matrix is non-singular, the  $A$  and  $A^{-1}$  will cancel. We also know that  $\|(F'(\vec{x}_n))^{-1}F(\vec{x}_n)\|$  is a good estimator for  $\|\vec{x}_n - \vec{x}^*\|$ , we can write the above equation as:

$$\|B^{-1}(B\vec{x}_{n+1} - B\vec{x}^*)\| \leq k\|B^{-1}(B\vec{x}_n - B\vec{x}^*)\|^2$$

where  $\vec{x}^*$  is the solution,  $k$  is a constant, and  $\vec{x}_n$ ,  $\vec{x}_{n+1}$  are the values of  $\vec{x}$  at steps  $n$  and  $n+1$ . At this point we note that the  $B$  and  $B^{-1}$  matrices cancel and any potential impact upon the convergence rate is nullified. From this analysis, Newton's method's convergence rate will not benefit from a linear

change of variables. An almost identical analysis of the Chord method yields the answer to potential performance increases.

The Chord method evaluates the Jacobian matrix once, prior the the iteration section of the code. The remainder of the code uses a set LU-decomposed version of the Jacobian to iterate to a solution. Since the Jacobian does not change over the life of the Chord method, only as small performance increase would be realized. The chord method is also governed by the equation:

$$\|\vec{x}_{n+1} - \vec{x}_n\| = \|(F'(\vec{x}_n))^{-1}F(\vec{x}_n)\|$$

Substituting in  $A$  and  $B$  into the above equation, we get:

$$\|\vec{x}_{n+1} - \vec{x}_n\| = \|(ABF'(B\vec{x}_0))^{-1}AF(B\vec{x}_n)\|$$

Since we know that the  $A$  matrix is non-singular, the  $A$  and  $A^{-1}$  will once again cancel. Using the fact that  $\|(F'(\vec{x}_0))^{-1}F(\vec{x}_n)\|$  is a good estimator for  $\|\vec{e}_n\|$   $\|\vec{e}_0\| = \|\vec{x}_n - \vec{x}^*\|$   $\|\vec{x}_0 - \vec{x}^*\|$ , we get:

$$\|B^{-1}(B\vec{x}_{n+1} - B\vec{x}^*)\| \leq k\|B^{-1}(B\vec{x}_0 - B\vec{x}^*)\| \|B^{-1}(B\vec{x}_n - B\vec{x}^*)\|$$

Once again, the non-singular  $B$  matrices cancel and have no impact upon the convergence rate of the Chord method.

To summarize, performing a linear change of variables on the Chord and Newton's methods has no impact upon the convergence rates of either methods. While there is a small possibility that performing a linear change of variables may improve the evaluation of the jacobian, it is likely that any performance gain would be offset by the cost of the matrix-vector multiplication required. Thus, it is not advisable to perform a linear change of variables.

## Question 2

*Write a program that solves single nonlinear equations with Newton's method, the chord method, and the secant method. For the the secant method, use  $x_{-1} = 0.99x_0$ . Apply your program to the following function/initial iterate combinations, document and explain your results:*

- (a)  $f(x) = 2x^2 - 5; x_0 = 10;$
- (b)  $f(x) = \sin(x) + x; x_0 = 0.5;$
- (c)  $f(x) = \cos(x); x_0 = 3$

The Newton's method, the Chord Method, and the Secant Method were all programmed in Matlab and used as nonlinear solvers for functions a,b, and c using the initial conditions given. Each of these codes can be found in the attached documentation, or listed in Appendix A below. All three methods used the following stopping criteria:

$$\|F(x)\| \leq \tau_r \|F(x_0)\| + \tau_a$$

where  $\tau_r$  is the relative tolerance compared to the initial norm of the function, and  $\tau_a$  is the absolute tolerance of the function. Both  $\tau_r$  and  $\tau_a$  were set to  $10^{-6}$  for this Question.

The first function was defined as:

$$f(x) = 2x^2 - 5 \quad x_0 = 10$$

This function is a basic parabola with a minima at  $x = 0$ , and  $F(x^*) = 0$  at  $x \approx \pm 1.5811$ . The three methods were all used to find the minima, and agree to down to the  $\approx 10^{-5}$  decimal place. The final  $x$  values of the three methods can be see in Figure 1. All three methods were able to converge to a common value for the first two functions. However, the periodicity of the cosine function in function c resulted in the failure of the Secant and Chord methods. Figure 2 shows the number of iterations required to converge for the three different methods. For Function a, Newton's method converged the quickest, but as can be seen by Figure 3, all three methods converged to the same point. The Chord method is a locally linearly convergent method and as a result, it took the longest to converge on the answer.

Function	Newton's	Chord Method	Secant Method
a	1.58113883	1.58113910	1.58113890
b	0.00000000	-0.00000000	0.00000000
c	-4.71238898	Did Not Converge	Did Not Converge

Figure 1: Final  $x$  values for the 3 functions

Function	Newton's	Chord Method	Secant Method
a	6	91	8
b	3	7	4
c	4	Did Not Converge	Did Not Converge

Figure 2: Number of Iterations required to converge

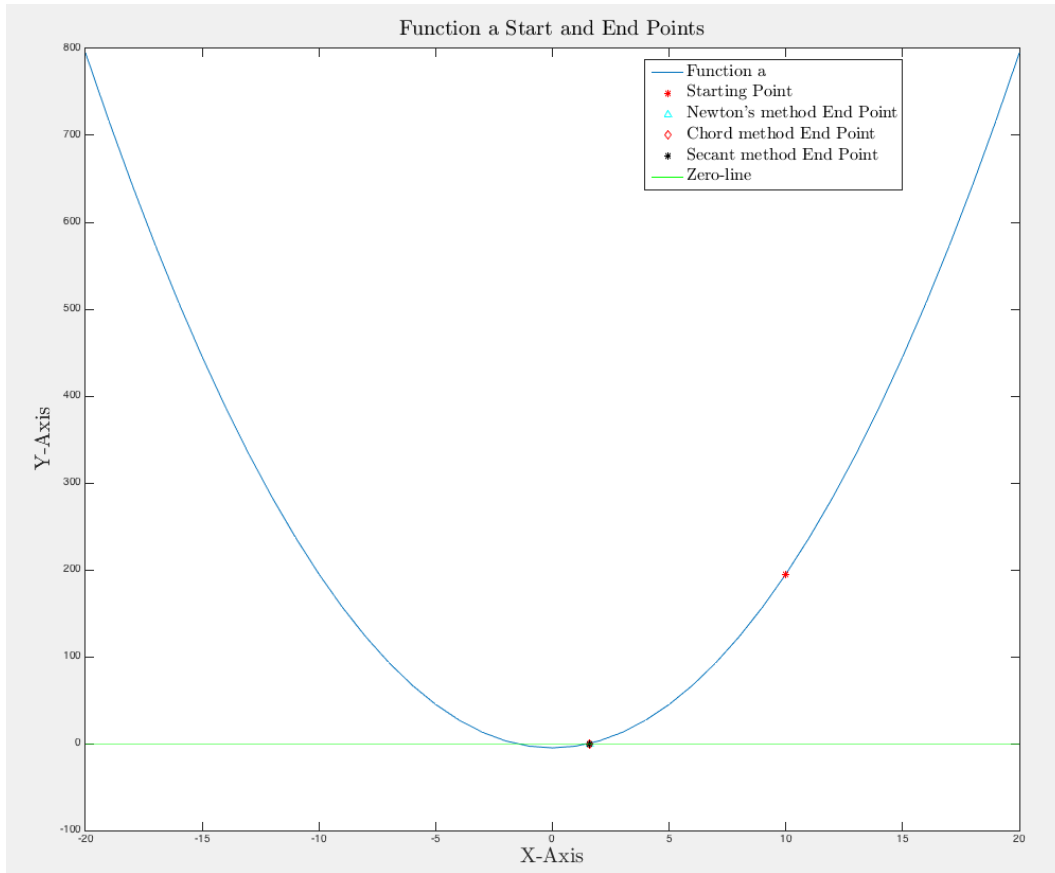


Figure 3: Initial conditions and final answer of Newton's Method, the Chord Method, and the Secant Method for Function a.

For Function b, once again, all three methods converged to the same point, as seen in Figure 4. As can be seen in Figure 2, Newton's method once again converges the fastest, but is only marginally quicker than the other two methods. For both Functions a and b, the starting point was close enough to the  $F(x^*) = 0$  point that the methods were able to converge.

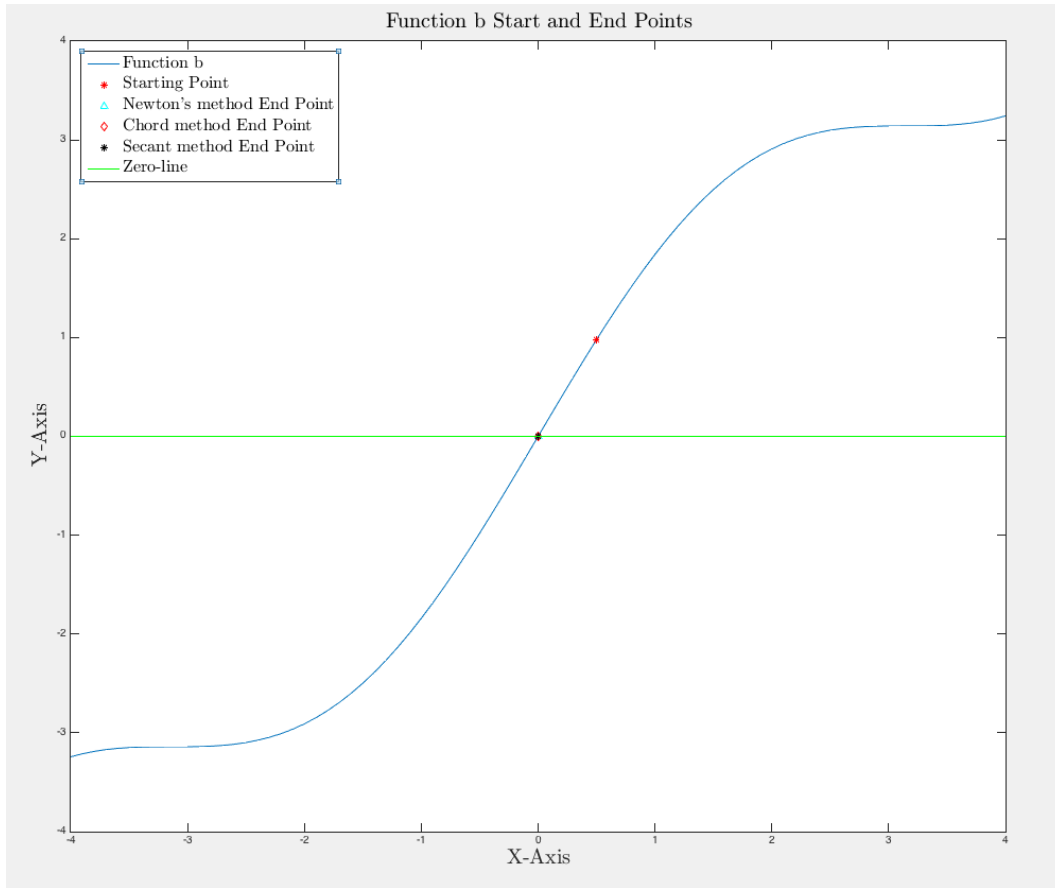


Figure 4: Initial conditions and final answer of Newton's Method, the Chord Method, and the Secant Method for Function b.

The Chord and Secant methods both had issues converging on Function c due to its periodic motion, as can be seen in Figures 1,2, and 5. While Newton's method did find a point where  $F(x) = 0$ , this was by no means the closest zero point. The Secant method diverged the fastest from the true solution, quickly reaching an  $x$  value equal to the computational precision of Matlab, which is  $2^{53}$ . At this point, the method is still working, but Matlab is unable to correctly calculate the next step. When  $a_n$  is being calculated, a divide-by-0 error occurs and the program exits. The Secant method works by approximating the derivative,  $F'(x)$  by using the previous step's  $x$  and  $F(x)$  values. However, with a periodic function, errors in the derivative can and will send the function off in an incorrect direction. If any  $x$  point is closest to to an extreme on the  $y$ -axis, the Secant — and Chord Method for that matter — will go careening off into the distance. This doesn't mean that they won't eventually converge, but if they do it will be at a  $F(x) = 0$  far from the start point.

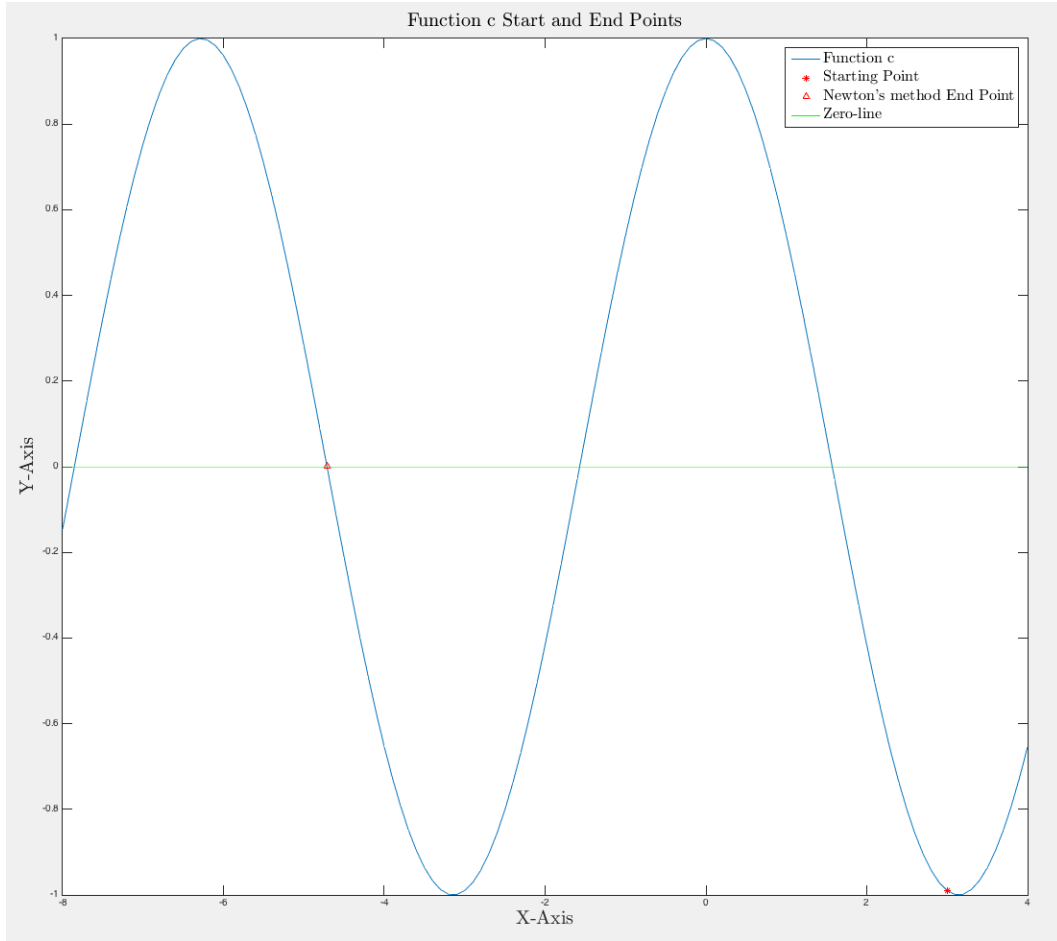


Figure 5: Initial conditions and final answer of Newton's Method, the Chord Method, and the Secant Method for Function c.

### Question 3

Assume that the standard assumptions hold, that the cost of a function evaluation is  $\mathcal{O}(N^2)$  floating-point operations, the cost of a Jacobian is  $\mathcal{O}(N)$  function evaluations, and that  $x_0$  is near enough to  $x$  so that the Newton iteration converges quadratically to  $x$ . Estimate what is the number of iteration needed to obtain  $\|e_n\| \leq \epsilon \|e_0\|$ , where  $\epsilon$  is a small tolerance value. What is the number of floating point operations required to get this accuracy?

Newton's Method converges quadratically to  $x$  and the error of the function is governed by the equation:

$$\|\vec{x}_{n+1} - \vec{x}^*\| \leq k \|\vec{x}_n - \vec{x}^*\|^2$$

where  $\vec{x}^*$  is the solution,  $k$  is a constant, and  $\vec{x}_n$ ,  $\vec{x}_{n+1}$  are the values of  $\vec{x}$  at steps  $n$  and  $n+1$ . Since  $\vec{x}^*$  is the solution,  $\vec{x}_n - \vec{x}^*$  can be written as:

$$\|\vec{e}_{n+1}\| \leq k \|\vec{e}_n\|^2 \leq (k \|\vec{e}_{n-1}\|^2)^2 \leq \dots$$

We know from these equations and the problem statement that the Newton iteration converges quadratically and that the norm of the error is dependent upon the previous step's error. Thus,

to get to  $\vec{x}_n$ , the error has converged by  $(k\|\vec{x}_0 - \vec{x}^*\|^2)^{n+1} = \|\vec{e}_n\|$ . The number of iterations necessary to get to the  $\epsilon$  tolerance is:

$$\frac{\|\vec{e}_n\|}{\|\vec{e}_0\|} = (k\|\vec{x}_0 - \vec{x}^*\|^2)^n = \epsilon$$

$$n = \frac{\log(\epsilon)}{\log(k\|\vec{e}_0\|^2)}$$

Assuming that evaluating the function,  $F(x)$  costs  $\mathcal{O}(N^2)$  FLOPS and evaluating the Jacobian is only  $\mathcal{O}(N)$  FLOPS, the number of of FLOPS required to get to an accuracy of  $\epsilon$  can be calculated. The breakdown of the FLOPS cost per step for Newton's method is given below:

$\vec{r} = \ F(\vec{x})\ $	$N^2 \text{ FLOPS}$
$F_x = \vec{r}$	$N \text{ FLOPS for an assignment}$
<i>While</i> $\ F_x\  > \tau_r \vec{r} + \tau_a$	$n \text{ iterations}$
$J = F'(\vec{x})$	$N \text{ FLOPS}$
$\vec{s} = -F(\vec{x})(F'(\vec{x}))^{-1}$	<i>Depends on how <math>\vec{s}</math> is solved. See Below</i>
$\vec{x} = \vec{x} + \vec{s}$	$N \text{ FLOPS}$
$F_x = F(\vec{x})$	$N^2 \text{ FLOPS}$
<i>end</i>	

The cost per iteration for Newton's method is highly dependent upon how  $\vec{s}$  is solved. From Homework 1, we know that the Conjugate Gradient method requires  $\approx \mathcal{O}(160N)$  FLOPS to reduce the error of the solver by a factor of  $10^{-3}$ . This value will be used as a rough estimate for the remained of this problem. Thus, adding up all of the FLOPS, it requires  $\mathcal{O}(N^2)$  FLOPS to set up the problem, and  $\mathcal{O}(N^2 + 162N)$  FLOPS per  $n$  iterations. To get to the  $\epsilon$  accuracy, it would take  $\mathcal{O}((\frac{\log(\epsilon)}{\log(k\|\vec{e}_0\|^2)})(N^2 + 162N) + N^2)$  FLOPS.

## Question 4

*Answer the questions in the previous problem for the chord method.*

Much of the previous problem's analysis carries over for this problem. One of the main differences between the Chord method and Newton's method is the Chord method converges at a linear rate in the area local to the solution. Because of this, the error equation presented in Problem 3 above changes to:

$$\frac{\|\vec{e}_n\|}{\|\vec{e}_0\|} = \frac{k\|\vec{e}_0\|\|\vec{e}_n\|}{k\|\vec{e}_0\|} = (\|\vec{x}_0 - \vec{x}^*\|)^n = \epsilon$$

$$n = \frac{\log(\epsilon)}{\log(k\|\vec{e}_0\|)}$$

Because the Chord method is linear, the quadratic term in the denominator changes into a linear term. Assuming that evaluating the function,  $F(x)$  costs  $\mathcal{O}(N^2)$  FLOPS and evaluating the Jacobian is only  $\mathcal{O}(N)$  FLOPS, the number of of FLOPS required to get to an accuracy of  $\epsilon$  can be calculated.

The breakdown of the FLOPS cost per step for the Chord method is given below:

$$\begin{array}{ll}
\vec{r} = ||F(\vec{x})|| & N^2 \text{ FLOPS} \\
F_x = \vec{r} & N \text{ FLOPS for an assignment} \\
J = F'(\vec{x}) & N \text{ FLOPS} \\
L, U \text{ decomposition of } J & \frac{N^3}{3} \text{ FLOPS for Cholesky factorization} \\
\text{While } ||F_x|| > \tau_r \vec{r} + \tau_a & n \text{ iterations} \\
L\vec{y} = -F(\vec{x}) & \text{Forward Substitution } N^2 \text{ FLOPS} \\
U\vec{s} = \vec{y} & \text{Back Substitution } N^2 \text{ FLOPS} \\
\vec{x} = \vec{x} + \vec{s} & N \text{ FLOPS} \\
F_x = F(\vec{x}) & N^2 \text{ FLOPS} \\
\text{end} & 
\end{array}$$

Once again, the cost per iteration is dependent on how  $\vec{s}$  is solved. Because the  $L, U$  decomposition has been calculated prior to the start of the iteration section, the cost is only incurred once. While  $\mathcal{O}(\frac{N^3}{3})$  is the worst-case for decomposing the Jacobian, in reality, this value is dependent on the sparsity of the matrix. If the Jacobian is very sparse, the cost for decomposing it will be much smaller. Inside the iteration section, the  $\vec{s}$  was solved using forward and backward substitution, but they could be easily calculated using two sequential Conjugate Gradient calls.

The total cost of using the Chord method with Cholesky factorization of the Jacobian and Forward/Backward substitution is  $\mathcal{O}((\frac{N^3}{3} + N^2 + 2N) + (\frac{\log(\epsilon)}{\log(k||\vec{e}_0||)})(3N^2 + N))$  FLOPS. This method, compared to Newton's method, has the potential for costing significantly more FLOPS to converge to a solution. However, if the Jacobian is extremely costly to evaluate, the Chord method has the advantage of only evaluating it once.



## Appendix A: Matlab code

### AutoRun

```
1 %% This is the auto running script for homework 2
2
3 atol=1e-8;
4 rtol=1e-8;
5 maxIt=10000;
6 %% Problem 2, part a
7
8 %Setup the problem
9 f1= @(x) 2*x.^2-5; % create a function handle for the problem
10 x0=10; %initialize the x val
11 xneg=0.99*x0; %for the secant method, x_(-1) val.
12 fprintf( '%s\n', 'Run 1 ');
13 [Nx1,newCount1]=Newton(f1,x0,atol,rtol);
14 [Cx1, Ccount1 ] = Chord( f1, x0,atol, rtol,maxIt);
15 [ Sx1, Scount1 ] = secant( f1, x0,xneg,atol, rtol,maxIt);
16 %% Problem 2, part b
17 fprintf( '%s\n', 'Run 2 ');
18
19 f2= @(x) sin(x)+x;
20 x0=0.5;
21 xneg=0.99*x0;
22
23 [Nx2,newCount2]=Newton(f2,x0,atol,rtol);
24 [Cx2, Ccount2 ] = Chord( f2, x0,atol, rtol,maxIt);
25 [ Sx2, Scount2 ] = secant( f2, x0,xneg,atol, rtol,maxIt);
26 %% Problem 2, part c
27
28 fprintf( '%s\n', 'Run 3 ');
29 f3= @(x) cos(x);
30 x0=0.3;
31 xneg=0.99*x0;
32
33
34 [Nx3,newCount3]=Newton(f3,x0,atol,rtol);
35 [Cx3, Ccount3 ] = Chord( f3, x0,atol, rtol,maxIt);
36 [ Sx3, Scount3 ] = secant( f3, x0,xneg,atol, rtol,maxIt);
37
38 %% Plotting section figure 1
39 figure;
40 x=-20:1:20;
41 plot(x,f1(x), 'DisplayName', 'Function a ');
42 hold on;
43 plot(10,f1(10), 'r*', 'DisplayName', 'Starting Point ');
44 plot(Nx1,f1(Nx1), 'c^', 'DisplayName', 'Newton's method End Point ');
45 plot(Cx1,f1(Cx1), 'rd', 'DisplayName', 'Chord method End Point ');
46 plot(Sx1,f1(Sx1), 'k*', 'DisplayName', 'Secant method End Point ');
47 plot(x,zeros(length(x),1), 'g', 'DisplayName', 'Zero-line ');
48
49 % Create ylabel
50 ylabel( { 'Y-Axis' }, 'FontSize',20, 'Interpreter', 'latex' );
51
52 % Create xlabel
53 xlabel( { 'X-Axis' }, 'FontSize',20, 'Interpreter', 'latex' );
```

```

54
55 % Create title
56 title({'Function a Start and End Points'}, 'FontSize', 20, ...
57       'Interpreter', 'latex');
58
59 % Create legend
60 legend1 = legend('show');
61 set(legend1, 'Interpreter', 'latex', 'FontSize', 16);
62
63 hold off;
64 %% Plotting section Figure 2
65 figure;
66 x=-4:.1:4;
67 plot(x, f2(x), 'DisplayName', 'Function b');
68
69 hold on;
70
71 plot(0.5, f2(0.5), 'r*', 'DisplayName', 'Starting Point');
72 plot(Nx2, f2(Nx2), 'c^', 'DisplayName', 'Newton's method End Point');
73 plot(Cx2, f2(Cx2), 'rd', 'DisplayName', 'Chord method End Point');
74 plot(Sx2, f2(Sx2), 'k*', 'DisplayName', 'Secant method End Point');
75
76 plot(x, zeros(length(x), 1), 'g', 'DisplayName', 'Zero-line');
77 % Create ylabel
78 ylabel({'Y-Axis'}, 'FontSize', 20, 'Interpreter', 'latex');
79
80 % Create xlabel
81 xlabel({'X-Axis'}, 'FontSize', 20, 'Interpreter', 'latex');
82
83 % Create title
84 title({'Function b Start and End Points'}, 'FontSize', 20, ...
85       'Interpreter', 'latex');
86
87 % Create legend
88 legend1 = legend('show');
89 set(legend1, 'Interpreter', 'latex', 'FontSize', 16);
90
91 hold off;
92 %% Plotting section Figure 3
93 figure;
94 x=-8:.1:4;
95 plot(x, f3(x), 'DisplayName', 'Function c');
96
97 hold on;
98
99 plot(03, f3(3), 'r*', 'DisplayName', 'Starting Point');
100 plot(Nx3, f3(Nx3), 'r^', 'DisplayName', 'Newton's method End Point');
101 % plot(Cx2, f3(Cx2), 'rd');
102 % plot(Sx3, f3(Sx3), 'k*');
103 plot(x, zeros(length(x), 1), 'g', 'DisplayName', 'Zero-line');
104 % Create ylabel
105 ylabel({'Y-Axis'}, 'FontSize', 20, 'Interpreter', 'latex');
106
107 % Create xlabel
108 xlabel({'X-Axis'}, 'FontSize', 20, 'Interpreter', 'latex');
109
110 % Create title
111 title({'Function c Start and End Points'}, 'FontSize', 20, ...
112       'Interpreter', 'latex');
113

```

```

114 % Create legend
115 legend1 = legend('show');
116 set(legend1, 'Interpreter', 'latex', 'FontSize', 16);
117
118 hold off;
119
120
121 %% Latex out
122 fprintf('function &Newton''s & Chord Method & Secant Method %s', '\\\\hline')
123 fprintf('a&%.8f& %.8f& %.8f%s \n', Nx1, Cx1, Sx1, '\\\\hline')
124 fprintf('b&%.8f& %.8f& %.8f%s \n', Nx2, Cx2, Sx2, '\\\\hline')
125 fprintf('c&%.8f& %s& %s%s \n', Nx3, 'Did Not Converge', 'Did Not Converge', '\\\\hline')
126 fprintf('\n');
127
128 fprintf('function &Newton''s & Chord Method & Secant Method %s\n', '\\\\hline')
129 fprintf('a&%d& %d& %d%s \n', newCount1, Ccount1, Scount1, '\\\\hline')
130 fprintf('b&%d& %d& %d%s \n', newCount2, Ccount2, Scount2, '\\\\hline')
131 fprintf('c&%d& %s& %s%s \n', newCount3, 'Did Not Converge', 'Did Not Converge', '\\\\hline')
132 fprintf('\n');

```

## Newton's Method

```

1 function [ x, numIts ] = Newton( fhandle, x0, atol, rtol)
2 %NEWTON Newton's method for non linear systems of equations
3 %%This method takes in the function handle to the system that needs to be
4 %%solved, the initial x value, and the tolerance.
5 %%Since the homework only requires a single nonlinear equation to be
6 %%solved, this method was not extended to cover a matrix.
7
8 r0=norm(fhandle(x0), inf);
9 x=x0;
10 fx=fhandle(x0);
11 numIts=0;
12 h=1e-5;
13
14 while norm(fx, inf)>rtol*r0+atol
15     numIts=numIts+1;
16     df=imag(fhandle(x+h*1i))/h;
17
18     %% Remove this for final runtime calcs!
19     % df2=(fhandle(x+h)-fx)/h;
20     % dabs=abs(df-df2);
21     % if(dabs>.001)
22     %     disp('ERROR WITH THE IMAGINARY STEP!!!');
23     % end
24
25     %% Solve for s. Since this is a 1-d problem, we can just divide by df.
26     %if this was a multi-dimensional matrix, we would need to use a linear
27     %solver to find s.
28
29     s=-fx/df;
30     x=x+s;
31     fx=fhandle(x);
32
33
34 end

```

```

35
36
37 end

```

## Chord Method

```

1 function [x, numIts] = Chord( fhandle, x0, atol, rtol, maxIt)
2 %CHORD The Chord method for non linear systems of equations
3 %%This method takes in the function handle to the system that needs to be
4 %%solved, the initial x value, and the tolerance.
5
6 r0=norm(fhandle(x0), inf);
7 x=x0;
8 fx=fhandle(x0);
9 numIts=0;
10 h=1e-5;
11
12 for i=1:length(x)
13     jacobian(i)=imag(fhandle(x+h*1i))/h;
14 end
15 [l,u]=lu(jacobian);
16
17
18 while norm(fx, inf)>rtol*r0+atol && numIts<maxIt
19     numIts=numIts+1;
20
21     if(mod(numIts,10000)==0)
22         fprintf('%s %d\n', 'Chord Method: Max number of iterations:', numIts);
23     end
24
25
26     %% Solve for s. Since this is a 1-d problem, we can just divide by df.
27     %%if this was a multi-dimensional matrix, we would need to use a linear
28     %%solver to find s.
29     y=-fx/l;
30     s=y/u;
31     x=x+s;
32     fx=fhandle(x);
33
34
35 end
36
37
38 end

```

## Secant Method

```

1 function [x, numIts] = secant( fhandle, x0, xneg, atol, rtol, maxIt)
2 %SECANT The secant method for non linear systems of equations
3 %%This method takes in the function handle to the system that needs to be
4 %%solved, the initial x value, and the tolerance.
5
6 r0=norm(fhandle(x0), inf);

```

```

7  x=x0;
8  x0=xneg;
9  fx=fhandle(x);
10 numIts=0;
11
12
13 while norm(fx,inf)>rtol*r0+atol && numIts<maxIt
14     numIts=numIts+1;
15
16     if (numIts==633)
17         disp(' ');
18     end
19
20     a=(fx - fhandle(x0))/(x - x0);
21     xn=x0;
22     x0=x;
23     x=x0 - fx/a;
24     fx=fhandle(x);
25
26
27 end
28
29
30 end

```