

Problem Sheet 4

Problem 1. Order of convergence from error recursion (core problem)

In [1, Exp. 2.3.26] we have observed *fractional* orders of convergence (\rightarrow [1, Def. 2.1.17]) for both the secant method, see [1, Code 2.3.25], and the quadratic inverse interpolation method. This is fairly typical for 2-point methods in 1D and arises from the underlying recursions for error bounds. The analysis is elaborated for the secant method in [1, Rem. 2.3.27], where a linearized error recursion is given in [1, Eq. (2.3.31)].

Now we suppose the recursive bound for the norms of the iteration errors

$$\|e^{(n+1)}\| \leq \|e^{(n)}\| \sqrt{\|e^{(n-1)}\|}, \quad (11)$$

where $e^{(n)} = x^{(n)} - x^*$ is the error of n -th iterate.

(1a) \square Guess the maximal order of convergence of the method from a numerical experiment conducted in MATLAB.

HINT:[1, Rem. 2.1.19]

Solution: See Listing 29.

Listing 29: MATLAB script for Sub-problem (1a)

```
1 % MATLAB test code for homework problem on "order of
  convergence from error
2 % recursion", outputs a table
3 e(1) = 1; e(2) = 0.8; % Any values < 1 possible, also
  random
```

```

4  for k=2:20, e(k+1) = e(k) * sqrt(e(k-1)); end
5  le = log(e); diff(le(2:end)) ./ diff(le(1:end-1)),

```

(1b) Find the maximal guaranteed order of convergence of this method through analytical considerations.

HINT: First of all note that we may assume equality in both the error recursion (11) and the bound $\|e^{(n+1)}\| \leq C\|e^{(n)}\|^p$ that defines convergence of order $p > 1$, because in both cases equality corresponds to a worst case scenario. Then plug the two equations into each other and obtain an equation of the type $\dots = 1$, where the left hand side involves an error norm that can become arbitrarily small. This implies a condition on p and allows to determine $C > 0$. A formal proof by induction (not required) can finally establish that these values provide a correct choice.

Solution: Suppose $\|e^{(n)}\| = C\|e^{(n-1)}\|^p$ (p is the largest convergence order and C is some constant).

Then

$$\|e^{(n+1)}\| = C\|e^{(n)}\|^p = C(C\|e^{(n-1)}\|^p)^p = C^{p+1}\|e^{(n-1)}\|^{p^2} \quad (12)$$

In (11) we may assume equality, because this is the worst case. Thus,

$$\|e^{(n+1)}\| = \|e^{(n)}\| \cdot \|e^{(n-1)}\|^{\frac{1}{2}} = C\|e^{(n-1)}\|^{p+\frac{1}{2}} \quad (13)$$

Combine (12) and (13),

$$C^{p+1}\|e^{(n-1)}\|^{p^2} = C\|e^{(n-1)}\|^{p+\frac{1}{2}}$$

i.e.

$$C^p\|e^{(n-1)}\|^{p^2-p-\frac{1}{2}} = 1. \quad (14)$$

Since (14) holds for each $n \geq 1$, we have

$$\begin{aligned} p^2 - p - \frac{1}{2} &= 0 \\ p &= \frac{1 + \sqrt{3}}{2} \quad \text{or} \quad p = \frac{1 - \sqrt{3}}{2} \quad (\text{dropped}). \end{aligned}$$

For C we find the maximal value 1.

Problem 2. Convergent Newton iteration (core problem)

As explained in [1, Section 2.3.2.1], the convergence of Newton's method in 1D may only be local. This problem investigates a particular setting, in which global convergence can be expected.

We recall the notion of a *convex function* and its geometric definition. A differentiable function $f : [a, b] \mapsto \mathbb{R}$ is convex, if and only if its graph lies on or above its tangent at any point. Equivalently, differentiable function $f : [a, b] \mapsto \mathbb{R}$ is convex, if and only if its derivative is non-decreasing.

Give a “graphical proof” of the following statement:

If $F(x)$ belongs to $C^2(\mathbb{R})$, is strictly increasing, is convex, and has a unique zero, then the Newton iteration [1, (2.3.4)] for $F(x) = 0$ is well defined and will converge to the zero of $F(x)$ for any initial guess $x^{(0)} \in \mathbb{R}$.

Solution: The sketches in Figure 8 discuss the different cases.

Problem 3. The order of convergence of an iterative scheme (core problem)

[1, Rem. 2.1.19] shows how to detect the order of convergence of an iterative method from a numerical experiment. In this problem we study the so-called [Steffensen's method](#), which is a derivative-free iterative method for finding zeros of functions in 1D.

Let $f : [a, b] \mapsto \mathbb{R}$ be twice continuously differentiable with $f(x^*) = 0$ and $f'(x^*) \neq 0$. Consider the iteration defined by

$$x^{(n+1)} = x^{(n)} - \frac{f(x^{(n)})}{g(x^{(n)})}, \quad \text{where} \quad g(x) = \frac{f(x + f(x)) - f(x)}{f(x)}.$$

(3a) \square Write a MATLAB script that computes the order of convergence to the point x^* of this iteration for the function $f(x) = xe^x - 1$ (see [1, Exp. 2.2.3]). Use $x^{(0)} = 1$.

Solution:

Listing 30: Order of convergence

```
1 clear all;
2 % different definitions of f with the same zero:
3 f = @(x) x.*exp(x)-1;
4 % f = @(x) exp(x)-1./x;           % f = @(x) x-exp(-x);
```

```

5 x0 = 1;
6 x_star = fzero(f,x0);
7
8 x = x0; upd = 1;
9 while (abs(upd) > eps)
10    fx = f(x(end));      % only 2 evaluations of f at each
11    step
12    if fx ~= 0;
13        upd = fx^2 / (f(x(end)+fx)-fx);
14        x = [x, x(end)-upd];
15    else upd = 0;
16    end
17 end
18 residual = f(x);
19 err      = abs(x-x_star);
20 log_err  = log(err);
21 ratios   = (log_err(3:end)-log_err(2:end-1))...
22     ./ (log_err(2:end-1)-log_err(1:end-2));
23 disp('x, err, residuals,ratios')
24 [x',err',residual',[0;0;rations']]
```

The output is

x	error e_n	$\frac{\log(e_{n+1}) - \log(e_n)}{\log(e_n) - \log(e_{n-1})}$
1.000000000000000	0.432856709590216	
0.923262600967822	0.356119310558038	
0.830705934728425	0.263562644318641	1.542345498206531
0.727518499997190	0.160375209587406	1.650553641703975
0.633710518522047	0.066567228112263	1.770024323911885
0.579846053882820	0.012702763473036	1.883754995643305
0.567633791946526	0.000490501536742	1.964598248590593
0.567144031581974	0.000000741172191	1.995899954235929
0.567143290411477	0.000000000001693	1.999927865685712
0.567143290409784	0.000000000000000	0.741551601040667

The convergence is obviously quadratic. As in the experiments shown in class, roundoff affects the estimated order, once the iteration error approaches the machine precision.

(3b) ☐ The function $g(x)$ contains a term like e^{xe^x} , thus it grows very fast in x and the method can not be started for a large $x^{(0)}$. How can you modify the function f (keeping the same zero) in order to allow the choice of a larger initial guess?

HINT: If f is a function and $h : [a, b] \rightarrow \mathbb{R}$ with $h(x) \neq 0, \forall x \in [a, b]$, then $(fh)(x) = 0 \Leftrightarrow f(x) = 0$.

Solution: The choice $\tilde{f}(x) = e^{-x}f(x) = x - e^{-x}$ prevents the blow up of the function g and allows to use a larger set of positive initial points. Of course, $\tilde{f}(x) = 0$ exactly when $f(x) = 0$.

Problem 4. Newton's method for $F(x) := \arctan x$

The merely local convergence of Newton's method is notorious, see [1, Section 2.4.2] and [1, Ex. 2.4.46]. The failure of the convergence is often caused by the overshooting of Newton correction. In this problem we try to understand the observations made in [1, Ex. 2.4.46].

(4a) ☐ Find an equation satisfied by the smallest positive initial guess $x^{(0)}$ for which Newton's method does not converge when it is applied to $F(x) = \arctan x$.

HINT: Find out when the Newton method oscillates between two values.

HINT: Graphical considerations may help you to find the solutions. See Figure 9: you should find an expression for the function g .

Solution: The function $\arctan(x)$ is **positive, increasing and concave** for positive x , therefore the first iterations of Newton's method with initial points $0 < x^{(0)} < y^{(0)}$ satisfy $y^{(1)} < x^{(1)} < 0$ (draw a sketch to see it). The function is **odd**, i.e., $\arctan(-x) = -\arctan(x)$ for every $x \in \mathbb{R}$, therefore the analogous holds for initial negative values ($y^{(0)} < x^{(0)} < 0$ gives $0 < x^{(1)} < y^{(1)}$). Moreover, opposite initial values give opposite iterations: if $y^{(0)} = -x^{(0)}$ then $y^{(n)} = -x^{(n)}$ for every $n \in \mathbb{N}$.

All these facts imply that, if $|x^{(1)}| < |x^{(0)}|$, then the absolute values of the following iterations will converge monotonically to zero. Vice versa, if $|x^{(1)}| > |x^{(0)}|$, then the absolute values of the Newton's iterations will diverge monotonically. Moreover, the iterations change sign at each step, i.e., $x^{(n)} \cdot x^{(n+1)} < 0$.

It follows that the smallest positive initial guess $x^{(0)}$ for which Newton's method does not converge satisfies $x^{(1)} = -x^{(0)}$. This can be written as

$$x^{(1)} = x^{(0)} - \frac{f(x^{(0)})}{f'(x^{(0)})} = x^{(0)} - (1 + (x^{(0)})^2) \arctan x^{(0)} = -x^{(0)}.$$

Therefore, $x^{(0)}$ is a zero of the function

$$g(x) = 2x - (1 + x^2) \arctan x \quad \text{with} \quad g'(x) = 1 - 2x \arctan x.$$

(4b) \square Use Newton's method to find an approximation of such $x^{(0)}$, and implement it with Matlab.

Solution: Newton's iteration to find the smallest positive initial guess reads

$$x^{(n+1)} = x^{(n)} - \frac{2x^{(n)} - (1 + (x^{(n)})^2) \arctan x^{(n)}}{1 - 2x^{(n)} \arctan x^{(n)}} = \frac{-x^{(n)} + (1 - (x^{(n)})^2) \arctan x^{(n)}}{1 - 2x^{(n)} \arctan x^{(n)}}.$$

The implementation in Matlab is given in Listing 31 (see also Figure 9).

Listing 31: Matlab Code for `NewtonArctan.m`

```

1  clear all;
2  x0 = 2;      % initial guess
3  r = 1;
4  while (r > eps)
5      x1 = x0 - ( 2*x0 - (1+x0^2) * atan(x0) ) /
             (1-2*x0*atan(x0));
6      % x1 = (-x0 + (1-x0^2) * atan(x0)) / (1-2*x0*atan(x0));
7      r = abs( (x1 - x0) / x1 );
8      x0 = x1;
9      fprintf( 'x0 = %16.14e , accuracy = %16.14e \n' ,
             x1, r );
10 end
11
12 figure;
13 x1 = x0-atan(x0)*(1+x0^2);    x2 = x1-atan(x1)*(1+x1^2);
14 X=[-2:0.01:2];
15 plot(X, atan(X), 'k',...
16       2*(X)-(1+(X).^2).*atan((X)), 'r--',...
17       [x0, x1, x1, x2, x2], [atan(x0), 0, atan(x1), 0,
             atan(x2)], ...
18       [x0,x1],[0,0], 'ro',[ -2,2], [0,0], 'k','linewidth',2);
19 legend('arctan', 'g', 'Newton critical iteration');axis
equal;
```

```
20 print -depsc2 'ex_NewtonArctan.eps'
```

Problem 5. Order- p convergent iterations

In [1, Section 2.1.1] we investigated the speed of convergence of iterative methods for the solution of a general non-linear problem $F(\mathbf{x}) = 0$ and introduced the notion of convergence of order $p \geq 1$, see [1, Def. 2.1.17]. This problem highlights the fact that for $p > 1$ convergence may not be guaranteed, even if the error norm estimate of [1, Def. 2.1.17] may hold for some $\mathbf{x}^* \in \mathbb{R}^n$ and all iterates $\mathbf{x}^{(k)} \in \mathbb{R}^n$.

Suppose that, given $\mathbf{x}^* \in \mathbb{R}^n$, a sequence $\mathbf{x}^{(k)}$ satisfies

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq C \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p \quad \forall k \quad \text{and some } p > 1.$$

(5a) Determine $\epsilon_0 > 0$ such that

$$\|\mathbf{x}^{(0)} - \mathbf{x}^*\| \leq \epsilon_0 \implies \lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*.$$

In other words, ϵ_0 tells us which distance of the initial guess from \mathbf{x}^* still guarantees local convergence.

Solution:

$$\lim_{k \rightarrow \infty} x^{(k)} = x^* \iff \lim_{k \rightarrow \infty} \|x^{(k)} - x^*\| = 0$$

Thus we seek an upper bound $B(k)$ for $\|x^{(k)} - x^*\|$ and claim that: $\lim_{k \rightarrow \infty} B(k) = 0$.

$$\begin{aligned} \|x^{(k)} - x^*\| &\leq C \|x^{(k-1)} - x^*\|^p \\ &\leq C \cdot C^p \|x^{(k-2)} - x^*\|^{p^2} \\ &\leq C \cdot C^p \cdot C^{p^2} \|x^{(k-3)} - x^*\|^{p^3} \\ &\vdots \\ &\leq C \cdots C^{p^{k-1}} \|x^{(0)} - x^*\|^{p^k} \\ &= C^{\sum_{i=0}^{k-1} p^i} \|x^{(0)} - x^*\|^{p^k} \\ &\stackrel{\text{geom. series}}{=} C^{\frac{p^k - 1}{p-1}} \|x^{(0)} - x^*\|^{p^k} \\ &\leq C^{\frac{p^k - 1}{p-1}} \epsilon_0^{p^k} = \underbrace{C^{\frac{1}{1-p}}}_{\text{const.}} \cdot \left(C^{\frac{1}{p-1}} \epsilon_0\right)^{p^k} = B(k) \end{aligned}$$

$$\lim_{k \rightarrow \infty} B(k) = 0 \iff C^{\frac{1}{p-1}} \epsilon_0 < 1$$

$$\implies 0 < \epsilon_0 < C^{\frac{1}{1-p}}$$

(5b) \square Provided that $\|\mathbf{x}^{(0)} - \mathbf{x}^*\| < \epsilon_0$ is satisfied, determine the minimal $k_{\min} = k_{\min}(\epsilon_0, C, p, \tau)$ such that

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| < \tau.$$

Solution: Using the previous upper bound and the condition τ , we obtain:

$$\|x^{(k)} - x^*\| \leq C^{\frac{1}{1-p}} \cdot \left(C^{\frac{1}{p-1}} \epsilon_0\right)^{p^k} < \tau$$

Solving for the minimal k (and calling the solution k_{\min}), with the additional requirement that $k \in \mathbb{N}$, we obtain:

$$\begin{aligned} \ln\left(C^{\frac{1}{1-p}}\right) + p^k \cdot \underbrace{\ln\left(C^{\frac{1}{p-1}} \epsilon_0\right)}_{<0} &< \ln \tau \\ k > \ln\left(\frac{\ln(\tau) + \frac{1}{p-1} \ln(C)}{\ln(C^{\frac{1}{p-1}} \epsilon_0)}\right) \cdot \frac{1}{\ln(p)}, \quad k_{\min} \in \mathbb{N} \\ k_{\min} &= \left\lceil \ln\left(\frac{\ln(\tau) + \frac{1}{p-1} \ln(C)}{\ln(C^{\frac{1}{p-1}} \epsilon_0)}\right) \cdot \frac{1}{\ln(p)} \right\rceil \end{aligned}$$

(5c) \square Write a MATLAB function

```
k_min = @(epsilon, C, p, tau) ...
```

and plot $k_{\min} = k_{\min}(\epsilon_0, \tau)$ for the values $p = 1.5, C = 2$. Test your implementation for every $(\epsilon_0, \tau) \in \text{linspace}(0, C^{\frac{1}{1-p}})^2 \cap (0, 1)^2 \cap \{(i, j) \mid i \geq j\}$

HINT: Use a MATLAB `pcolor` plot and the commands `linspace` and `meshgrid`.

Solution: See `k_min_plot.m`.

Listing 32: Matlab Code for `k_min_plot.m`, Problem 5.

<code>1</code> C <code>2</code> p	<code>= 2;</code> <code>= 1.5;</code>
--------------------------------------	--

```

3  eps_max          = C^(1/(1-p));
4
5  ngp              = 100; % number of grid points
6  eps_lin          = linspace(0, eps_max, ngp);
7  tau_lin          = linspace(0, eps_max, ngp);
8  [eps_msh,tau_msh] =
    meshgrid(eps_lin(2:(end-1)),tau_lin(2:(end-1)));
9
10 kmin = @(eps, C, p, tau) ceil(log( ( log(tau) +
    (1/(p-1)).*log(C) ) ./ log(C^(1/(p-1)).*eps) ) ./
    log(p) );
11 k = kmin(eps_msh, C, p, tau_msh);
12
13 % Consider only gridpoints where: eps larger as tau
14 for ne = 1:ngp-2
15     for nt = 1:ngp-2
16         if (ne > nt)
17             k(ne,nt) = 0;
18         end
19     end
20 end
21
22 % Plotting
23 pcolor(eps_msh,tau_msh,k)
24 colorbar()
25 title('Minimal number of iterations for error < \tau')
26 xlabel('\epsilon_0')
27 ylabel('\tau')
28 xlim([0,eps_max])
29 ylim([0,eps_max])
30 shading flat

```

Problem 6. Code quiz

A frequently encountered drudgery in scientific computing is the use and modification of poorly documented code. This makes it necessary to understand the ideas behind the code first. Now we practice this in the case of a simple iterative method.

(6a) ☐ What is the purpose of the following MATLAB code?

```

1   function y = myfn(x)
2     log2 = 0.693147180559945;
3
4     y = 0;
5     while (x > sqrt(2)), x = x/2; y = y + log2; end
6     while (x < 1/sqrt(2)), x = x*2; y = y - log2; end
7     z = x-1;
8     dz = x*exp(-z)-1;
9     while (abs(dz/z) > eps)
10       z = z+dz;
11       dz = x*exp(-z)-1;
12     end
13     y = y+z+dz;

```

Solution: The MATLAB code computes $y = \log(x)$, for a given x . The program can be regarded as Newton iterations for finding the zero of

$$f(z) = e^z - x \quad (15)$$

where z is unknown and x is given.

(6b) ☐ Explain the rationale behind the two `while` loops in lines #5, 6.

Solution: The purpose of the two while loops is to shift the function values of (15) and modify the initial $z_0 = x - 1$ in such a way that good convergence is reached (according to the function derivative).

(6c) ☐ Explain the loop body of lines #10, 11.

Solution: The `while`-loop computes the zero of (15) using Newton iterations

$$\begin{cases} z^{(0)} = x - 1, \\ z^{(n+1)} = z^{(n)} - \frac{e^{z^{(n)}} - x}{e^{z^{(n)}}}, \quad n \geq 1 \end{cases}$$

(6d) ☐ Explain the conditional expression in line #9.

Solution: This is a correction based termination criterium, see [1, § 2.1.26].

(6e) Replace the `while`-loop of lines #9 through #12 with a fixed number of iterations that, nevertheless, guarantee that the result has a relative accuracy `eps`.

Solution: Denote the zero of $f(z)$ with z^* , and $e^{(n)} = z^{(n)} - z^*$. Use Taylor expansion of $f(z)$, $f'(z)$:

$$\begin{aligned} f(z^{(n)}) &= e^{z^*} e^{(n)} + \frac{1}{2} e^{z^*} (e^{(n)})^2 + O((e^{(n)})^3) \\ &= xe^{(n)} + \frac{1}{2} x (e^{(n)})^2 + O((e^{(n)})^3), \\ f'(z^{(n)}) &= e^{z^*} + e^{z^*} e^{(n)} + O((e^{(n)})^2) \\ &= x + xe^{(n)} + O((e^{(n)})^2) \end{aligned}$$

Considering Newton's Iteration $z^{(n+1)} = z^{(n)} - \frac{f(z^{(n)})}{f'(z^{(n)})}$,

$$\begin{aligned} e^{(n+1)} &= e^{(n)} - \frac{f(z^{(n)})}{f'(z^{(n)})} \\ &= e^{(n)} - \frac{xe^{(n)} + \frac{1}{2}x(e^{(n)})^2 + O((e^{(n)})^3)}{x + xe^{(n)} + O((e^{(n)})^2)} \\ &\doteq \frac{1}{2}(e^{(n)})^2 \\ &= \dots \\ &\doteq \left(\frac{1}{2}\right)^{1+2+\dots+2^n} (e^0)^{2^{n+1}} \\ &= 2 \cdot \left(\frac{1}{2}\right)^{2^{n+1}} (e^0)^{2^{n+1}} \\ &= 2 \cdot \left(\frac{1}{2}e^0\right)^{2^{n+1}}, \end{aligned}$$

where $e^0 = z^0 - z^* = x - 1 - \log(x)$. So it is enough for us to determine the number n of iteration steps by $|\frac{e^{(n+1)}}{\log x}| = \text{eps}$. Thus

$$\begin{aligned} n &= \frac{\log(\log(2) - \log(\log(x)) - \log(\text{eps})) - \log(\log(2) - \log(|e^0|))}{\log(2)} - 1 \\ &\approx \frac{\log(-\log(\text{eps})) - \log(-\log(|e^0|))}{\log(2)} - 1 \end{aligned}$$

The following code is for your reference.

```

function y = myfn(x)
    log2 = 0.693147180559945;
    y = 0;
    while (x > sqrt(2)), x = x/2; y = y + log2; end
    while (x < 1/sqrt(2)), x = x*2; y = y - log2; end
    z = x-1;
    dz = x*exp(-z)-1;
    e0=z-log(x);
    k=(log(-log(eps))-log(-log(abs(e0))))/log(2);
    for i=1:k
        z = z+dz;
        dz = x*exp(-z)-1;
    end
    y = y+z+dz;

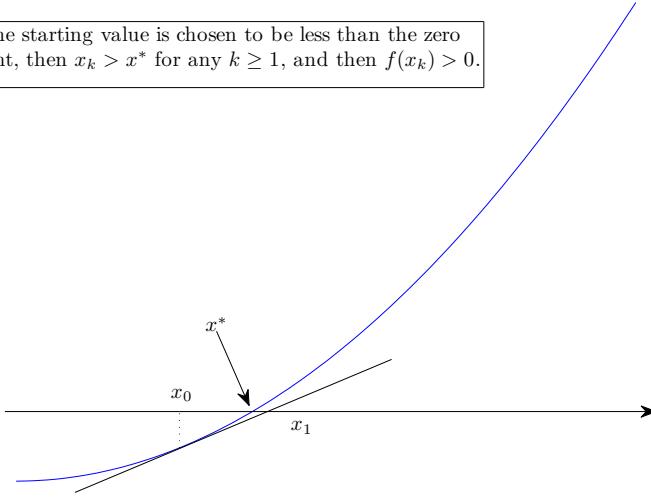
```

Issue date: 08.10.2015

Hand-in: 15.10.2015 (in the boxes in front of HG G 53/54).

Version compiled on: October 15, 2015 (v. 1.0).

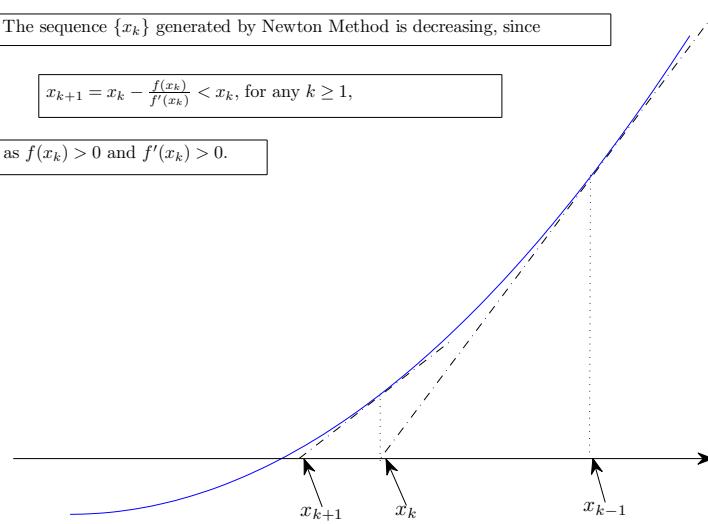
If the starting value is chosen to be less than the zero point, then $x_k > x^*$ for any $k \geq 1$, and then $f(x_k) > 0$.



The sequence $\{x_k\}$ generated by Newton Method is decreasing, since

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} < x_k, \text{ for any } k \geq 1,$$

as $f(x_k) > 0$ and $f'(x_k) > 0$.



Every tangent of the convex function
 $g_k(x) = f(x_k) + f'(x_k)(x - x_k)$ is
below the function graph $f(x)$. Thus $x_k > x^*$ for every $k \geq 1$ and
then the limit of the sequence exists with $\lim_{k \rightarrow \infty} x_k = y^* \geq x^*$. If
 $y^* > x^*$, then $g_k(y^*) = f(y^*) > f(x^*) = 0$ which means that the
method can not be stopped at this moment (as the stopping
criterion $|x_{k+1} - x_k| = \frac{f(x_k)}{f'(x_k)}$ is not small enough). Thus $y^* = x^*$.

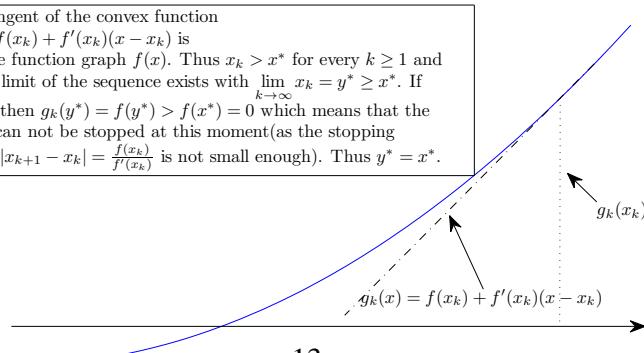


Figure 8: Graphical proof for Problem 2.

Figure 9: Newton iterations with $F(x) = \arctan(x)$ for the critical initial value $x^{(0)}$

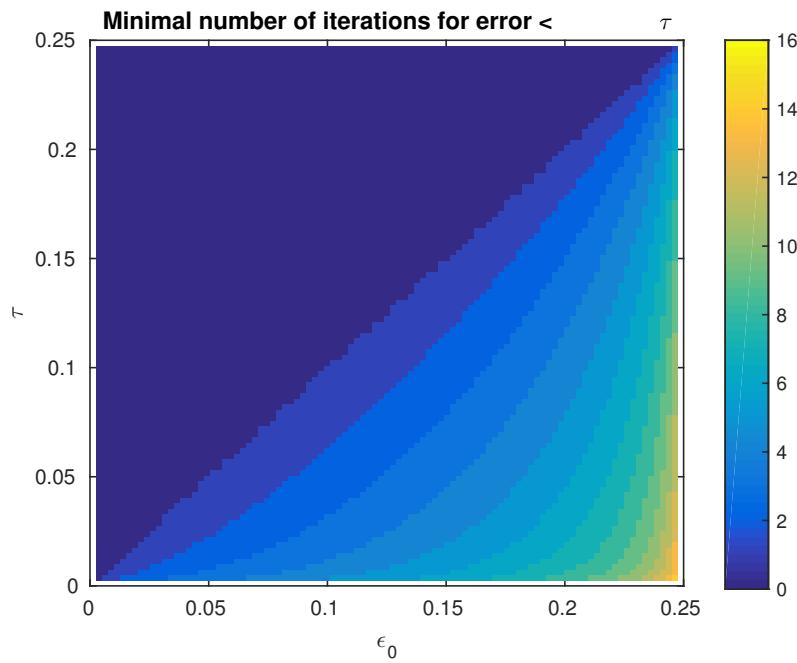
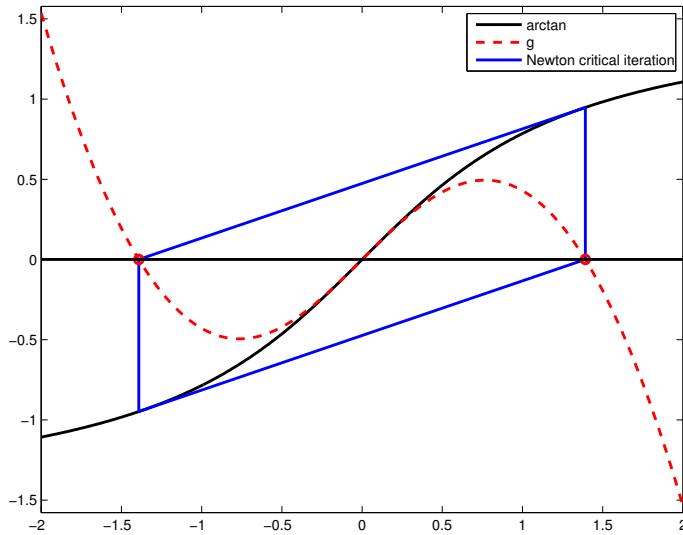


Figure 10: Minimum number of iterations given τ and ϵ_0 , Problem 5.