

SOLVING THE EQUATION $f(x) = 0$

In this lecture we present the bisection algorithm in its general setting and use it to prove two theorems: Bolzano's Theorem and the Intermediate Value Theorem. We then present another algorithm, the fixed point iteration, and use it to prove the Contraction Mapping Theorem.

Note: the figures are missing, I did not have time to do them yet.

1. THE BISECTION ALGORITHM

We consider real-valued functions of a real variable, $f : \mathbf{R} \rightarrow \mathbf{R}$, or $f : [a, b] \rightarrow \mathbf{R}$ where $[a, b]$ is an interval of real numbers.

We first consider equations of the form $f(x) = 0$.

Theorem. (*Bolzano's theorem*) Assume that $f : [a, b] \rightarrow \mathbf{R}$ is Lipschitz continuous and that $f(a)f(b) < 0$ (opposite signs). Then there is a real number $\bar{x} \in [a, b]$ such that $f(\bar{x}) = 0$. If f is strictly monotone (increasing or decreasing) then \bar{x} is unique (the only such number).

The proof is a *constructive proof*, in contrast to, for example, a proof by contradiction. This means that the proof describes how the number \bar{x} is constructed. It is organized in the following four steps:

- (1) an algorithm which produces an approximating sequence $\{x_i\}$;
- (2) a proof that $\{x_i\}$ is a Cauchy sequence so that we get a real number (decimal expansion) $\bar{x} = \lim_{i \rightarrow \infty} x_i$.
- (3) a proof that \bar{x} solves the equation: $\{f(x_i)\}$ is also a Cauchy sequence and $f(\bar{x}) = \lim_{i \rightarrow \infty} f(x_i) = 0$;
- (4) a proof that \bar{x} is unique (in the case of strictly monotone function).

Remember these steps, we will use the same kind of proof many times. The first three steps give the *existence* of a solution. The last step gives *uniqueness* of the solution of the equation $f(x) = 0$. One important consequence of uniqueness is that all approximating sequences will converge to the same solution \bar{x} . That is, the solution is independent of the choice of algorithm or approximating sequence (independent of the construction).

Proof. Step 1. We use the bisection algorithm with starting points a and b .

Step 2. We obtain two sequences $\{x_i\}_{i=0}^{\infty}$ and $\{X_i\}_{i=0}^{\infty}$ with

- (1) $|x_i - X_i| \leq (b - a)2^{-i}$,
- (2) $|x_i - x_j| \leq (b - a)2^{-i}, \quad j > i$,
- (3) $|X_i - X_j| \leq (b - a)2^{-i}, \quad j > i$.

The inequalities (2) and (3) mean that x_i and X_i are Cauchy sequences: $|x_i - x_j| \rightarrow 0, |X_i - X_j| \rightarrow 0$ as $i, j \rightarrow \infty$, and we get decimal expansions (real numbers)

$$\bar{x} = \lim_{i \rightarrow \infty} x_i, \quad \bar{X} = \lim_{i \rightarrow \infty} X_i.$$

The inequality (1) means that the decimal expansions are the same:

$$\bar{x} = \lim_{i \rightarrow \infty} x_i = \lim_{i \rightarrow \infty} X_i = \bar{X}.$$

Step 3. The Lipschitz continuity of f gives

$$|f(x_i) - f(x_j)| \leq L|x_i - x_j| \leq L(b - a)2^{-i}, \quad j > i,$$

which means that $f(x_i)$ is a Cauchy sequence, $|f(x_i) - f(x_j)| \rightarrow 0$ as $i, j \rightarrow \infty$. It gives a decimal expansion (real number) which we denote $f(\bar{x})$:

$$f(\bar{x}) = \lim_{i \rightarrow \infty} f(x_i).$$

We must show that this is equal to 0. To do this we note that the distance between $f(x_i)$ and 0 is less than or equal to the distance between $f(x_i)$ and $f(X_i)$. This is because $f(x_i)$ and $f(X_i)$ have opposite signs. For example, if $f(x_i) < 0$ and $f(X_i) > 0$:

$$|f(x_i) - 0| = -f(x_i) < -f(x_i) + f(X_i) = |f(x_i) - f(X_i)|.$$

In any case:

$$|f(x_i) - 0| \leq |f(x_i) - f(X_i)| \leq L|x_i - X_i| \leq L(b - a)2^{-i} \rightarrow 0,$$

where we also used the Lipschitz condition and (1). This means that

$$\lim_{i \rightarrow \infty} f(x_i) = 0,$$

or in other words $f(\bar{x}) = 0$.

Step 4. Assume now that f is strictly increasing. (The case of decreasing function can be handled in a similar way.) This means that $x < y$ implies $f(x) < f(y)$. Assume also that we have two different solutions, i.e., $\bar{x}_1 < \bar{x}_2$ with $f(\bar{x}_1) = f(\bar{x}_2) = 0$. But this is a contradiction to the strict monotonicity. Therefore, there is only one solution. \square

FIGURE 1. Non-monotone function with three roots.

FIGURE 2. Monotone function with unique root.

We now consider equations of the form $f(x) = y$.

Theorem. (*The Intermediate Value Theorem*) *If $f : [a, b] \rightarrow \mathbf{R}$ is Lipschitz continuous and the real number y is between $f(a)$ and $f(b)$, then there is a real number $\bar{x} \in [a, b]$ such that $f(\bar{x}) = y$. If f is strictly monotone, then \bar{x} is unique.*

The theorem says that a Lipschitz continuous function takes all values between its end-point values. A discontinuous function may skip some value.

Proof. If $f(a) = f(b)$ then $y = f(a) = f(b)$ and we can choose $\bar{x} = a$ or $\bar{x} = b$. Otherwise we apply Bolzano's theorem to the function $F(x) = f(x) - y$. Clearly, $F : [a, b] \rightarrow \mathbf{R}$ is Lipschitz with the same constant as f , and $F(a)F(b) < 0$ because y is between $f(a)$ and $f(b)$. Also, F is strictly monotone if f is strictly monotone. The conclusion now follows from Bolzano's theorem: there is \bar{x} such that $F(\bar{x}) = f(\bar{x}) - y = 0$, and it is unique if f is strictly monotone. \square

FIGURE 3. Equation $f(x) = y$ with multiple solutions.

FIGURE 4. Equation $f(x) = y$ with unique solution.

FIGURE 5. The discontinuous function skips the value y .

2. INVERSE FUNCTION

Assume now that the function in the Intermediate Value Theorem is strictly increasing. (Decreasing functions can be discussed in the same way.) Let us write $A = f(a)$, $B = f(b)$. We note the following consequences of the theorem.

The function takes all values y in the interval $[A, B]$ and it takes no values outside $[A, B]$. This means that the range of the function is exactly $R_f = [A, B]$.

The equation $f(x) = y$ has a unique solution x for any $y \in [a, b]$. Since x is unique, this defines a function $y \mapsto x$. More precisely, we can define a function

$$g : [A, B] \rightarrow \mathbf{R}$$

$$x = g(y), \quad \text{where } x \text{ is the unique solution of } f(x) = y.$$

(Remember that a function must have a *unique* value for *each* element of the domain of definition.) We then say that f is invertible (“inverterbar”) and the function g is called the the inverse of f . It is denoted $g = f^{-1}$,

$$f^{-1} : [A, B] \rightarrow \mathbf{R}$$

$$x = f^{-1}(y), \quad \text{where } x \text{ is the unique solution of } f(x) = y.$$

Note that

$$D_{f^{-1}} = R_f = [A, B], \quad R_{f^{-1}} = D_f = [a, b],$$

$$f(f^{-1}(y)) = y, \quad y \in [A, B]; \quad f^{-1}(f(x)) = x, \quad x \in [a, b].$$

Of course, we may interchange the roles of x and y and write $y = f^{-1}(x)$, $x \in [A, B]$.

Warning: f^{-1} is pronounced “f inverse”, and it is not the same as “f to the minus one” $(f)^{-1} = 1/f$. It is confusing that “inverse” and “raised to -1 ” are written in the same way, but it is a tradition in mathematics, and we have to live with it and be careful when we use it.

We have now proved the following.

Theorem. *A strictly monotone function is invertible.*

Example. The function $f : [-2, 2] \rightarrow \mathbf{R}$, $f(x) = x^2$, is not invertible because the equation $x^2 = y$ has two solutions $x = \pm\sqrt{y}$. In the next section we shall see that f becomes invertible if we restrict its domain of definition to the nonnegative numbers.

3. THE FUNCTION \sqrt{x}

We introduce the square root function and investigate its properties.

The function $f : [0, b] \rightarrow \mathbf{R}$, $f(x) = x^2$, is strictly increasing because $0 \leq x < y$ implies $x - y < 0$ and $x + y > 0$ so that $x^2 - y^2 = (x - y)(x + y) < 0$. Hence the equation $x^2 = y$ has a unique solution $x = \sqrt{y}$ for any $y \in [0, b^2]$. Therefore, f is invertible with $f^{-1}(y) = \sqrt{y}$.

This can be done for any b , the uniqueness implies that the bisection algorithm gives the same result no matter what starting points $a = 0$ and $b > 0$ we use. So the function $f : \mathbf{R}_+ \rightarrow \mathbf{R}$, $f(x) = x^2$, is invertible with inverse $f^{-1} : \mathbf{R}_+ \rightarrow \mathbf{R}$, $f^{-1}(x) = \sqrt{x}$. Note: $D_{\sqrt{}} = \mathbf{R}_+$ and $R_{\sqrt{}} = \mathbf{R}_+$.

Warning: $f^{-1}(x) = \sqrt{x}$ is not the same as $(f)^{-1}(x) = (f(x))^{-2} = x^{-2} = 1/x^2$ although they are written in the same way.

FIGURE 6. The function $y = \sqrt{x}$.

Since $x^2 = a$ and $y^2 = b$ implies $(xy)^2 = ab$ and $(x/y)^2 = a/b$, we conclude

$$\sqrt{ab} = \sqrt{a}\sqrt{b}, \quad \sqrt{\frac{a}{b}} = \frac{\sqrt{a}}{\sqrt{b}}.$$

We now investigate the Lipschitz continuity of \sqrt{x} :

$$\begin{aligned} |\sqrt{x} - \sqrt{y}| &= \left\{ \text{multiply by the "conjugate expression" } \sqrt{x} + \sqrt{y} \right\} \\ &= \left| \frac{(\sqrt{x} - \sqrt{y})(\sqrt{x} + \sqrt{y})}{\sqrt{x} + \sqrt{y}} \right| = \left| \frac{x - y}{\sqrt{x} + \sqrt{y}} \right| \quad \left\{ \text{by the conjugate rule} \right\} \\ &= \frac{1}{\sqrt{x} + \sqrt{y}} |x - y| \leq \frac{1}{\sqrt{\delta} + \sqrt{\delta}} |x - y| = \frac{1}{2\sqrt{\delta}} |x - y| \quad \text{for } x, y \geq \delta. \end{aligned}$$

We conclude that the square root function is Lipschitz continuous with constant $L = 1/(2\sqrt{\delta})$ on any interval of the form $[\delta, \infty)$ with $\delta > 0$. That is, on any interval that stays way from 0. But it is not Lipschitz on its whole domain of definition $\mathbf{R}_+ = [0, \infty)$. This is clear from the previous calculation, but is also seen in the graph where the slope is infinite at 0.

4. THE POWER FUNCTION x^r FOR $r \in \mathbf{Q}$

The power function ("potensfunktioner") x^r with rational exponent $r = p/q$ is defined by uniquely solving the equation $y^q = x^p$ for $x \geq 0$. This is possible because the function $f: \mathbf{R}_+ \rightarrow \mathbf{R}$, $f(x) = x^p$, is strictly increasing. Can you prove this? Hint:

$$x^p - y^p = (x - y)(x^{p-1} + x^{p-2}y + \dots + xy^{p-2} + y^{p-1}).$$

The unique solution is denoted $y = x^{p/q} = x^r$. It satisfies

$$x^r x^s = x^{r+s}, \quad \frac{x^r}{x^s} = x^{r-s}, \quad r, s \in \mathbf{Q}.$$

We shall define x^r with real exponent $r \in \mathbf{R}$ later in ALA-b.

5. THE FIXED POINT ITERATION

An algebraic equation can be written in two equivalent ways (meaning that they have the same solutions).

1. $f(x) = 0$. A solution \bar{x} is called a *root* of f or a *zero* of f ("nollställe till f ").

Example. The function $f(x) = x^2 - 2$ has two roots $\bar{x}_1 = \sqrt{2}$, $\bar{x}_2 = -\sqrt{2}$.

Algorithm: For equations of this form we have the bisection algorithm.

2. $x = g(x)$. A solution \bar{x} is called a *fixed point* ("fixpunkt") of g . The equation is called a fixed point equation.

Example. The function $g(x) = 2/x$ has two fixed points, $\bar{x}_1 = \sqrt{2}$, $\bar{x}_2 = -\sqrt{2}$, because $\pm\sqrt{2} = \pm \frac{\sqrt{2}\sqrt{2}}{\sqrt{2}} = \pm\sqrt{2}$.

Algorithm: a natural algorithm for a fixed point equation is to choose a starting point x_0 and then compute x_i according to $x_i = g(x_{i-1})$. This is called the fixed point iteration. We hope that the sequence x_i converges to a fixed point. This works sometimes and sometimes not.

Example. With $g(x) = x/2 + 1/x$ and $x_0 = 1$ we get $x_1 = g(x_0) = 3/2$, $x_2 = 15/12$ and so on. This is easy to try in MATLAB:

```
>> format long
>> x=1
>> x=x/2+1/x
>> x=x/2+1/x
>> x=x/2+1/x
>> x=x/2+1/x
```

Do this now!! Does it converge? Do you recognize a decimal expansion?

Example. With $g(x) = 2/x$ and $x_0 = 1$ we get $x_1 = g(x_0) = 2$, $x_2 = 1$, $x_3 = 2$, i.e., we get the sequence $\{1, 2, 1, 2, \dots\}$ which is divergent.

FIGURE 7. The functions $y = g(x)$ and $y = x$.

Note that the equations $x^2 - 2 = 0$, $x = x/2 + 1/x$, and $x = 2/x$ are equivalent (multiply the last two equations by x to see this). We can rewrite equations between the two forms in many ways. For example, $x = g(x)$ can be written as $x - g(x) = 0$. On the other hand, $f(x) = 0$ can be written $x = x + \alpha f(x)$ where α or $\alpha(x)$ is any nonzero number or function. The challenge is to find a “good” choice of α so that the fixed point iteration is convergent. Later, when we discuss Newton’s method, we shall find an optimal choice of α . In fact, $x = x/2 + 1/x$ is obtained from $x^2 - 2 = 0$ by using $\alpha(x) = -1/(2x)$ and we will learn why this is a good choice.

So when does the fixed point iteration work? It turns out that one important condition is that the Lipschitz constant $L = L_g$ of g is strictly smaller than 1. More precisely, we shall assume that g is Lipschitz on an interval I with constant $L < 1$, i.e.,

$$|g(x) - g(y)| \leq L|x - y| \quad \text{for } x, y \in I \text{ and with } L < 1.$$

Such a function is called a *contraction mapping* (or just contraction). Note that the distance between the images $g(x)$ and $g(y)$ is smaller than the distance between x and y .

Example. For $g(x) = 2/x$ we have

$$|g(x) - g(y)| = 2 \left| \frac{y - x}{xy} \right| = \frac{2}{|x||y|} |x - y| \leq \frac{1}{2} |x - y|, \quad \forall x, y \geq 2,$$

so that g is a contraction with $L = 1/2$ on $I = [2, \infty)$. Note that we used that $z = xy \in [4, \infty)$ so that $\frac{2}{|x||y|} = \frac{2}{z} \in [0, \frac{1}{2}]$.

Example. For $g(x) = x/2 + 1/x$ we have

$$|g(x) - g(y)| = \left| \frac{x - y}{2} + \frac{y - x}{xy} \right| = \left| \frac{1}{2} - \frac{1}{xy} \right| |x - y|.$$

Now let us consider $x, y \in [1, 2]$ for example. Write $z = xy$. Then $z \in [1, 4]$ and we get $\frac{1}{2} - \frac{1}{xy} = \frac{1}{2} - \frac{1}{z} \in [-\frac{1}{2}, \frac{1}{4}]$ with absolute value $|\frac{1}{2} - \frac{1}{xy}| \leq \frac{1}{2}$. Therefore

$$|g(x) - g(y)| = \left| \frac{1}{2} - \frac{1}{xy} \right| |x - y| \leq \frac{1}{2} |x - y|, \quad \forall x, y \in [1, 2],$$

so that g is a contraction with $L = 1/2$ on $I = [1, 2]$.

6. THE CONTRACTION MAPPING THEOREM

Remember that the bisection algorithm leads to Bolzano’s theorem. The fixed point iteration also leads to a theorem.

Remember that a closed interval I is an interval that contains its endpoints (if it has any). A closed interval can be of the following types:

$$I = [a, b] \quad (\text{closed and bounded interval})$$

$$I = [a, \infty), I = (-\infty, b], I = (-\infty, \infty) = \mathbf{R}, \quad (\text{closed and unbounded intervals})$$

Theorem. (*The Contraction Mapping Theorem*) Assume that I is a closed interval and that $g : I \rightarrow I$ is a contraction mapping. Then g has a unique fixed point $\bar{x} \in I$. The fixed point is obtained as the limit of the fixed point iteration, $x_i = g(x_{i-1})$, for any starting point $x_0 \in I$.

It is important that the target set I is the same as the domain of definition, $g : I \rightarrow I$; it guarantees that the sequence does not jump out of the interval I where g is a contraction. It is also important that I is closed; it guarantees that $\bar{x} = \lim x_i \in I$.

Proof. The proof follows the four steps of a constructive proof that we mentioned before.

Step 1. An algorithm: we use the fixed point iteration. Take an arbitrary point $x_0 \in I$ and compute $x_i = g(x_{i-1})$.

Step 2. A proof that $\{x_i\}$ is a Cauchy sequence. We must estimate $|x_i - x_j|$ for $j > i$. Consider first the distance between two consecutive elements of the sequence:

$$|x_{k+1} - x_k| = |g(x_k) - g(x_{k-1})| \leq L|x_k - x_{k-1}|.$$

Here we used the fact that the x_k stay in I and g is a contraction on I . Therefore

$$|x_{k+1} - x_k| \leq L|x_k - x_{k-1}|.$$

Since $L < 1$ this means that x_{k+1}, x_k are closer to each other than x_k, x_{k-1} . In the same way:

$$|x_k - x_{k-1}| \leq L|x_{k-1} - x_{k-2}|.$$

By repeating this we get

$$\begin{aligned} |x_{k+1} - x_k| &\leq L|x_k - x_{k-1}| \\ &\leq L^2|x_{k-1} - x_{k-2}| \\ &\leq L^3|x_{k-2} - x_{k-3}| \\ &\leq \dots \leq L^k|x_1 - x_0|, \end{aligned}$$

that is

$$(4) \quad |x_{k+1} - x_k| \leq L^k|x_1 - x_0|.$$

Now consider $|x_i - x_j|$ for $j > i$. We have

$$\begin{aligned} x_i - x_j &= x_i - x_{i+1} + x_{i+1} - x_{i+2} + x_{i+2} - \dots - x_{j-2} + x_{j-2} - x_{j-1} + x_{j-1} - x_j \\ &= \sum_{k=i}^{j-1} (x_k - x_{k+1}). \end{aligned}$$

Such a sum is called a telescope sum because all terms cancel except the first and the last. Applying the triangle inequality to the sum and using (4) we have

$$|x_i - x_j| \leq |x_i - x_{i+1}| + \dots + |x_{j-1} - x_j| = \sum_{k=i}^{j-1} |x_k - x_{k+1}| \leq |x_1 - x_0| \sum_{k=i}^{j-1} L^k.$$

This is a geometric sum given by the well-known formula:

$$\sum_{k=i}^{j-1} L^k = L^i(1 + L + \dots + L^{j-i-1}) = L^i \frac{1 - L^{j-i}}{1 - L}.$$

Therefore

$$(5) \quad |x_i - x_j| \leq |x_1 - x_0| L^i \frac{1 - L^{j-i}}{1 - L} \leq |x_1 - x_0| L^i \frac{1}{1 - L},$$

because $0 \leq 1 - L^{j-i} \leq 1$ for $j > i$. Since $L < 1$ we have $L^i \rightarrow 0$ and hence $|x_i - x_j| \rightarrow 0$ as $i \rightarrow \infty$ with $j > i$. Thus x_i is a Cauchy sequence and we get a decimal expansion (real number)

$$\bar{x} = \lim_{i \rightarrow \infty} x_i,$$

which belongs to I because I is closed.

Step 3. Proof that \bar{x} is a fixed point. We have

$$|g(x_i) - g(x_j)| \leq L|x_i - x_j| \rightarrow 0, \quad i, j \rightarrow \infty,$$

so that $g(x_i)$ is a Cauchy sequence. We get a real number which we denote $g(\bar{x})$:

$$g(\bar{x}) = \lim_{i \rightarrow \infty} g(x_i).$$

We must show that $\lim_{i \rightarrow \infty} g(x_i) = \bar{x}$ so that $g(\bar{x}) = \bar{x}$. But

$$|\bar{x} - g(x_i)| = |\bar{x} - x_{i+1}| \rightarrow 0, \quad i \rightarrow \infty.$$

This means that $\lim_{i \rightarrow \infty} g(x_i) = \bar{x}$ and hence $g(\bar{x}) = \bar{x}$.

Step 4. Uniqueness. Assume that we have two fixed points $\bar{x}_1, \bar{x}_2 \in I$. Then

$$|\bar{x}_1 - \bar{x}_2| = |g(\bar{x}_1) - g(\bar{x}_2)| \leq L|\bar{x}_1 - \bar{x}_2|.$$

which implies

$$(1 - L)|\bar{x}_1 - \bar{x}_2| \leq 0.$$

But $1 - L > 0$ so the only possibility is $|\bar{x}_1 - \bar{x}_2| = 0$. In other words: $\bar{x}_1 = \bar{x}_2$. So there is only one fixed point in I .

Note that the uniqueness of \bar{x} implies that we get the same limit no matter which starting point x_0 we choose. \square

Example. We have seen that $g(x) = 2/x$ is a contraction on the closed interval $I = [2, \infty)$. But $x_0 = 3$ gives $x_1 = 2/3 \notin I$ so the sequence jumps out. The sequence jumps back and forth between 3 and $2/3$; it does not converge.

Example. We have seen that $g(x) = x/2 + 1/x$ is a contraction with $L = 1/2$ on $I = [1, 2]$. We check that $g : I \rightarrow I$. If $x \in I = [1, 2]$, i.e., $1 \leq x \leq 2$, then $x/2 \leq 1$ and $1/x \leq 1$ so that $x/2 + 1/x \leq 1 + 1 = 2$. Also $x/2 \geq 1/2$ and $1/x \geq 1/2$ so that $x/2 + 1/x \geq 1$. Therefore $g(x) \in I = [1, 2]$. The contraction mapping theorem says that g has a unique fixed point in $I = [1, 2]$. What is it?

7. WHEN DO WE STOP THE ITERATION?

We stop the iteration when the distance between two consecutive iterates is less than a given tolerance, $|x_i - x_{i+1}| \leq \text{tol}$. Then we expect that a certain number of decimals have been fixed in the decimal expansion \bar{x} . For example, with $|x_i - x_{i+1}| \leq 10^{-N-1}$ we expect approximately N decimals to be fixed.

This is justified by the calculation (with $j > i$ as usual)

$$\begin{aligned} |x_i - x_j| &\leq \sum_{k=i}^{j-1} |x_k - x_{k+1}| \leq |x_i - x_{i+1}| \sum_{k=i}^{j-1} L^{k-i} \\ &\leq \frac{1 - L^{j-i}}{1 - L} |x_i - x_{i+1}| \leq \frac{1}{1 - L} |x_i - x_{i+1}|, \end{aligned}$$

which is done in the same way as (5) but using $|x_{k+1} - x_k| \leq L^{k-i}|x_i - x_{i+1}|$ instead of (4). Therefore

$$|x_i - x_j| \leq \frac{1}{1 - L} |x_i - x_{i+1}| \leq \frac{1}{1 - L} \text{tol}, \quad j > i.$$

or by letting $j \rightarrow \infty$

$$|x_i - \bar{x}| \leq \frac{1}{1 - L} |x_i - x_{i+1}| \leq \frac{1}{1 - L} \text{tol}.$$

So the number of fixed decimals after i steps is determined by $\frac{1}{1-L} \text{tol}$. Note that this number is bigger than tol , much bigger if L is near 1, so we get fewer decimals than tol itself indicates. The reason for this is that we are looking at the *residual* $x_i - x_{i+1} = x_i - g(x_i)$ which measures how well x_i satisfies the equation $x - g(x) = 0$. The size of the residual is then magnified by the factor $\frac{1}{1-L}$ when we use it to estimate the error in x_i .

8. HOW FAST IS THE CONVERGENCE?

From the construction of the sequence x_i we have

$$|x_i - \bar{x}| = |g(x_{i-1}) - g(\bar{x})| \leq L|x_{i-1} - \bar{x}|.$$

Therefore the error is reduced by the factor $L < 1$ in each step. The smaller L is, the faster the convergence. This is called *linear convergence*. The bisection algorithm also converges linearly: the error is reduced by a factor $1/2$ in each step.

Example. For $g(x) = x - x^2/2 + 1$ on $I = [1, 3/2]$ we have $L = 1/2$. This follows from

$$g(x) - g(y) = (1 - \frac{1}{2}(x + y))(x - y).$$

Here $2 \leq x + y \leq 3$ so that $-\frac{1}{2} \leq 1 - \frac{1}{2}(x + y) \leq 0$ with absolute value $|1 - \frac{1}{2}(x + y)| \leq 1/2$.

If we compute a few iterations in MATLAB with $x_0 = 1$ and for each iteration compute $|x_i - \sqrt{2}|$ we see that the error is reduced approximately by a factor 1/2 in each step. This is what I got. The first column is x_i and the second is $|x_i - \sqrt{2}|$.

1.000000000000000	0.41421356237310
1.500000000000000	0.08578643762690
1.375000000000000	0.03921356237310
1.429687500000000	0.01547393762690
1.40768432617188	0.00652923620122
1.41689674509689	0.00268318272380
1.41309855196381	0.00111501040929
1.41467479318270	0.00046123080961
1.41402240794944	0.00019115442365
1.41429272285787	0.00007916048478

Note that $\sqrt{2}$ is computed by the MATLAB function `sqrt(2)` which is also an approximation but with approximately 16 correct decimals. So we can use it to test the accuracy of our computation.

The convergence is sometimes much faster than this. As the example $g(x) = x/2 + 1/x$ shows. Compute and check this!! This is what I got: again the first column is x_i and the second is $|x_i - \sqrt{2}|$.

1.000000000000000	0.41421356237310
1.500000000000000	0.08578643762690
1.416666666666667	0.00245310429357
1.41421568627451	0.00000212390141
1.41421356237469	0.00000000000159
1.41421356237309	0.00000000000000
1.41421356237309	0.00000000000000
1.41421356237309	0.00000000000000
1.41421356237309	0.00000000000000
1.41421356237309	0.00000000000000
1.41421356237309	0.00000000000000

In fact it is possible to show that

$$|x_i - \bar{x}| \leq K|x_i - \bar{x}|^2,$$

where K is some number. This is called *quadratic convergence* and means that the error is reduced by a factor which is proportional to the error itself. So the speed of converges increases as x_i approaches \bar{x} . We will learn later why it is so fast. See Chapter Newton's method.

9. ADVANTAGES AND DISADVANTAGES

A disadvantage with the fixed point iteration is that it is often difficult to find a suitable interval I where the iteration converges. This is usually very easy for the bisection algorithm: just plot the function and pick two points where the function has opposite signs.

There are two advantages:

- (1) the fixed point iteration can be very fast. The bisection algorithm always converges linearly but not faster than that.
- (2) the fixed point iteration also works for systems of equations. This is not true for bisection.

We shall return to these advantages later.

/stig