

CONVERGENCE THEORY

* FIXED-POINT ITERATIONS $x_{n+1} = g(x_n)$

Using the same solution R , we write an expression for the error in the n -th iteration

$$R - x_{n+1} = R - g(x_n) = g(R) - g(x_n)$$

Multiplying and dividing by $(R - x_n)$

$$R - x_{n+1} = \frac{g(R) - g(x_n)}{R - x_n} (R - x_n)$$

Using the mean-value theorem

$$R - x_{n+1} = g'(z_n) \cdot (R - x_n), \text{ where } z_n \in [x_n, R]$$

Now, denoting $e_n = R - x_n$, we have

$$|e_{n+1}| = |g'(z_n)| \cdot |e_n|$$

Remarks

- * if $|g'(z_n)| < 1$ on some interval, the fixed-point iterations will converge for initial values in that interval (This explains why certain rearrangements of $f(x)=0$ may result in divergent iterations)
- * the error at a given iteration is a fraction of the error at the previous iteration \Rightarrow linear convergence
- * to ensure convergence, the function $g(x)$ must be contractive, i.e., Lipschitz-continuous with the constant $L < 1$

* Newton's Method

Newton's method uses iterations similar to the fixed-point approach

$$x_{n+1} = x_n - \frac{g(x_n)}{g'(x_n)} = g(x_n)$$

Iterates converge if $|g'(x)| < 1$

Or we

$$|g'(x)| = \left| \frac{g(x) \cdot g''(x)}{[g'(x)]^2} \right| < 1 \quad // \text{ condition for convergence}$$

Error relation: $R - x_{n+1} = g(R) - g(x_n)$

Taylor-expand $g(x_n)$ about $x=R$ up to the second order

$$\begin{cases} g(x_n) = g(R) + g'(R)(x_n - R) + \frac{g''(s)}{2}(x_n - R)^2 \\ g(R) = 0 \Rightarrow g'(R) = \frac{g(R) \cdot g''(R)}{[g'(R)]^2} = 0 \end{cases} \quad s \in [x_n, R]$$

$$g(x_n) = g(R) + \frac{g''(s)}{2} (x_n - R)^2$$

Putting this into the error relation

$$|e_{n+1}| = |g(R) - g(x_n)| = \left| \frac{g''(s)}{2} \right| |e_n|^2$$

Remarks

* smallness of the second derivative required for convergence

* The error at a given iteration is the square of the error at the previous iteration — quadratic convergence

SECANT AND FALSE POSITION METHODS

In both cases the consecutive iterates can be written as

$$x_{n+1} = x_n - \frac{g(x_n)}{g(x_n) - g(x_{n-1})} (x_n - x_{n-1}) = \underline{g(x_n, x_{n-1})}$$

Following a similar approach as above, but with more tedious details, we obtain

$$x_{n+1} = \frac{g(s_1, s_2)}{2} x_{n-1}, \quad s_1, s_2 \in [x_n, R]$$

Error is proportional to the sum of previous errors.

Remark

- * It can be shown that the rate of convergence is faster than linear, but slower than quadratic
 \Rightarrow Superlinear convergence

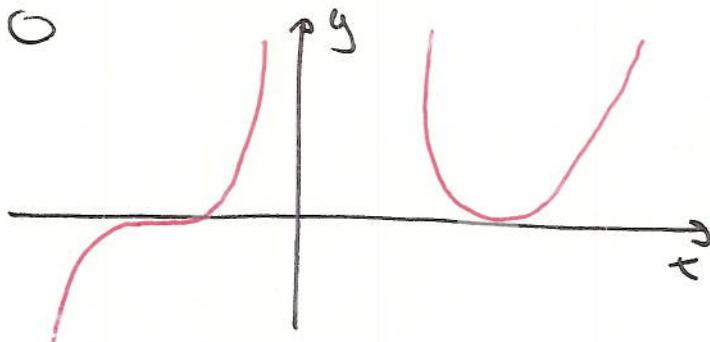
General Remarks

- * Higher-order methods (Newton's method) are faster, but come with more stringent requirements on the function $f(x)$ (boundedness of higher derivatives)
- * One can combine different methods:
 - Start with a slower, but more robust method
 - Switch to a faster method close to the root.

* problems occur when multiple roots are present, i.e.,
when $g(R) = g'(R) = \dots = 0$

In Newton's method

$$g'(x) = \frac{g(x) \cdot g''(x)}{[g'(x)]^2} \sim \frac{0}{0}$$



doesn't vanish and second-order convergence is lost.

1.3 Numerical Linear Algebra - A Short Review

(MATH 2703, Grasselli & Belinsonsky Sections 2.3-2.6)

Problem: We are interested in the solution of problems of the type

$$Ax = b, \quad x, b \in \mathbb{R}^N, \quad A \in \mathbb{R}^{N \times N}$$

Assumption

N-large

General remarks about existence of solutions:

- * if $\text{Det}(A) \neq 0$ a unique solution x exists
 $\{ \text{if } b=0, x=0$
- * if $\text{Det}(A)=0$, let $x = x_0 + x_1$,
 - $Ax_0 = 0$ has m solutions defined by m multiplicative constants, $m =$ the number of distinct eigenvalues
 - $Ax_1 = b$ - additional solutions exist if $b \perp N(A^T)$ - Frobenius alternative

(Cramer's Rule - Sol'n involving determinants (ugly),
but useful for large N)

* Elimination Methods

Start with the system below; The goal is to convert it to an upper-triangular form

$$\begin{array}{l} 4x_1 - 2x_2 + x_3 = 15 \quad | \cdot 3 \rightarrow \\ -3x_1 - x_2 + 4x_3 = 8 \quad | \cdot 4 \rightarrow \\ x_1 - x_2 + 3x_3 = 13 \quad | \cdot 4 \end{array} \Rightarrow$$

$$\Rightarrow \begin{array}{l} 4x_1 - 2x_2 - x_3 = 15 \\ -10x_2 + 19x_3 = 77 \\ -2x_2 + 11x_3 = 37 \end{array} \begin{array}{l} | \cdot 2 \rightarrow \\ | \cdot (-10) \rightarrow \end{array} \Rightarrow \left\{ \begin{array}{l} 4x_1 - 2x_2 + x_3 = 15 \\ -10x_2 + 19x_3 = 77 \\ -72x_3 = -216 \end{array} \right.$$

Given the triangular form, the solution can be found trivially via back-substitution: $x_3 = 3, x_2 = -2, x_1 = 3$

The same goal can be achieved by performing elementary row operations on the corresponding system matrix:

- multiplication of a row by a constant
- adding a multiple of a row to another row
- interchanging the orders of two rows

Problem

- the magnitude of the coefficients grows; this may become an issue for large systems (round-off errors)