

## Chapter 7:

We are solving the ODE

$$\vec{y}'(t) = \vec{f}(t, \vec{y}(t)), \quad t \geq 0, \quad \vec{y}(t_0) = \vec{y}_0.$$

Assume  $\vec{f}$  is non-linear.

The backward Euler method is

$$\vec{y}_{n+1} = \vec{y}_n + h \vec{f}(t_{n+1}, \vec{y}_{n+1}).$$

The trapezoidal rule is

$$\vec{y}_{n+1} = \vec{y}_n + \frac{1}{2} h (\vec{f}(t_{n+1}, \vec{y}_{n+1}) + \vec{f}(t_n, \vec{y}_n)).$$

The general form of these two methods can be written as

$$\vec{\omega} = h \vec{g}(\vec{\omega}) + \vec{\beta}, \quad \vec{\omega} \in \mathbb{R}^d \quad \text{--- (*)}$$

For the backward Euler method, we have

$$\vec{g}(\vec{\omega}) = \vec{f}(t_{n+1}, \vec{\omega}), \quad \vec{\beta} = \vec{y}_n.$$

For the trapezoidal rule, we have

$$\vec{g}(\vec{\omega}) = \frac{1}{2} \vec{f}(t_{n+1}, \vec{\omega}), \quad \vec{\beta} = \vec{y}_n + \frac{1}{2} h \vec{f}(t_n, \vec{y}_n).$$

Implicit Runge-Kutta methods and multi-step methods may also be written in this form.

Goal: Given some initial guess  $\vec{\omega}^{[0]}$ , find algorithm

$$\vec{\omega}^{[i+1]} = \vec{s}(\vec{\omega}^{[i]}), \quad i = 0, 1, 2, \dots$$

such that:

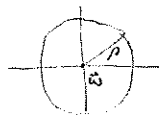
- 1)  $\vec{\omega}^{[i]} \rightarrow \vec{\omega}$  which is solution to (\*);
- 2) cost is small;
- 3) algorithm converges fast.

The most simple is fixed-point iteration (also known as functional iteration sometimes):

$$\vec{\omega}^{[i+1]} = h \vec{g}(\vec{\omega}^{[i]}) + \vec{\beta}, \quad i = 0, 1, 2, \dots$$

Notation:  $\|\cdot\|$  is a vector norm and  $B_\rho(\vec{\omega})$  is a closed ball of radius  $\rho > 0$  centered at  $\vec{\omega}$ :

$$B_\rho(\vec{\omega}) = \{ \vec{u} \in \mathbb{R}^d; \|\vec{u} - \vec{\omega}\| \leq \rho \}.$$

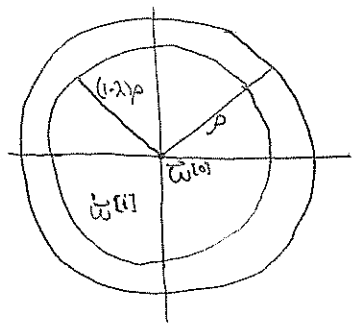


Theorem 7.1: Let  $h > 0$ ,  $\bar{\omega}^{[0]} \in \mathbb{R}^d$ , and suppose there are  $\lambda \in (0, 1)$  and  $\rho > 0$  such that:

- (i)  $\|\bar{g}(\bar{v}) - \bar{g}(\bar{u})\| \leq \frac{1}{h} \|\bar{v} - \bar{u}\| \quad \forall \bar{v}, \bar{u} \in \mathcal{B}_\rho(\bar{\omega}^{[0]})$ ;
- (ii)  $\bar{\omega}^{[1]} \in \mathcal{B}_{(1-\lambda)\rho}(\bar{\omega}^{[0]})$ .

Then:

- (a)  $\bar{\omega}^{[i]} \in \mathcal{B}_\rho(\bar{\omega}^{[0]})$  (solution does not diverge);
- (b)  $\hat{\omega} = \lim_{i \rightarrow \infty} \bar{\omega}^{[i]}$  exists and obeys  $\bar{\omega} = h\bar{g}(\bar{\omega}) + \bar{\beta}$ ;
- (c)  $\hat{\omega}$  is unique.



Proof: We prove first that

$$\|\bar{\omega}^{[i+1]} - \bar{\omega}^{[i]}\| \leq \lambda^i (1-\lambda)\rho.$$

This holds for  $i=0$  by (ii). We proceed by induction. Assume the statement is true  $\forall m=0, 1, \dots, i-1$ . Then

$$\begin{aligned} \|\bar{\omega}^{[i+1]} - \bar{\omega}^{[i]}\| &= \left\| (h\bar{g}(\bar{\omega}^{[i]}) + \bar{\beta}) - (h\bar{g}(\bar{\omega}^{[i-1]}) + \bar{\beta}) \right\| \\ &= h \|\bar{g}(\bar{\omega}^{[i]}) - \bar{g}(\bar{\omega}^{[i-1]})\| \\ &\leq \lambda \|\bar{\omega}^{[i]} - \bar{\omega}^{[i-1]}\| \\ &\leq \lambda \lambda^{i-1} (1-\lambda)\rho \\ &= \lambda^i (1-\lambda)\rho. \end{aligned}$$

(by (i))  
 $\uparrow$   
 don't know (i) is applicable

This completes the induction.

Note that

$$\bar{\omega}^{[i+1]} - \bar{\omega}^{[0]} = \sum_{j=0}^i (\bar{\omega}^{[j+1]} - \bar{\omega}^{[j]}), \quad i=0, 1, 2, \dots$$

and so

$$\begin{aligned} \|\bar{\omega}^{[i+1]} - \bar{\omega}^{[0]}\| &= \left\| \sum_{j=0}^i (\bar{\omega}^{[j+1]} - \bar{\omega}^{[j]}) \right\| \\ &\leq \sum_{j=0}^i \|\bar{\omega}^{[j+1]} - \bar{\omega}^{[j]}\| \end{aligned}$$

this should all be one single induction

$$\begin{aligned}
&\leq \sum_{j=0}^i \lambda^j (1-\lambda) \rho \\
&= (1-\lambda^{i+1}) \rho \\
&\leq \rho.
\end{aligned}$$

This proves (a).

Now we turn to (b). Note

$$\bar{\omega}^{[i+k]} - \bar{\omega}^{[i]} = \sum_{j=0}^{k-1} (\bar{\omega}^{[i+j+1]} - \bar{\omega}^{[i+j]})$$

and so

$$\begin{aligned}
\|\bar{\omega}^{[i+k]} - \bar{\omega}^{[i]}\| &\leq \sum_{j=0}^{k-1} \|\bar{\omega}^{[i+j+1]} - \bar{\omega}^{[i+j]}\| \\
&\leq \sum_{j=0}^{k-1} \lambda^{i+j} (1-\lambda) \rho \\
&= \lambda^i (1-\lambda^k) \rho.
\end{aligned}$$

bad wording

Since  $\lambda \in (0, 1)$ , there exists  $k$  large enough such that for any  $\varepsilon > 0$

$$\|\bar{\omega}^{[i+k]} - \bar{\omega}^{[i]}\| < \varepsilon.$$

Therefore  $(\bar{\omega}^{[i]})_{i=0}^{\infty}$  is a Cauchy sequence. Since  $\mathcal{B}_\rho(\bar{\omega}^{[0]})$  is complete, there is  $\hat{\omega} \in \mathcal{B}_\rho(\bar{\omega}^{[0]})$  such that  $\|\hat{\omega} - \bar{\omega}^{[i]}\| \rightarrow 0$  as  $i \rightarrow \infty$ . This proves (b).

We prove (c) by contradiction. Assume  $\bar{\omega}^* \in \mathcal{B}_\rho(\bar{\omega}^{[0]})$  with  $\bar{\omega}^* \neq \hat{\omega}$  and  $\bar{\omega}^* = h\bar{g}(\bar{\omega}^*) + \bar{\beta}$ . Then  $\|\bar{\omega}^* - \hat{\omega}\| > 0$  implies

$$\begin{aligned}
\|\bar{\omega}^* - \hat{\omega}\| &= \|(h\bar{g}(\bar{\omega}^*) + \bar{\beta}) - (h\bar{g}(\hat{\omega}) + \bar{\beta})\| \\
&= h \|\bar{g}(\bar{\omega}^*) - \bar{g}(\hat{\omega})\| \\
&\leq \lambda \|\bar{\omega}^* - \hat{\omega}\| \\
&< \|\bar{\omega}^* - \hat{\omega}\|.
\end{aligned}$$

■

2016 03 31

Let's look at condition (i) of theorem 7.1. Let's assume that  $\bar{g}$  is smoothly differentiable. Then there is  $\tau \in (0, 1)$  such that

$$\bar{g}(\bar{v}) - \bar{g}(\bar{u}) = \frac{\partial \bar{g}(\tau\bar{v} + (1-\tau)\bar{u})}{\partial \bar{\omega}} (\bar{v} - \bar{u}).$$

Hence

$$\|\bar{g}(\bar{v}) - \bar{g}(\bar{u})\| \leq \left\| \frac{\partial \bar{g}(\tau\bar{v} + (1-\tau)\bar{u})}{\partial \bar{\omega}} \right\| \|\bar{v} - \bar{u}\|.$$

So (i) says something about the magnitude of  $h$  in relation to

$$\left\| \frac{\partial \vec{g}}{\partial \vec{\omega}} \right\|.$$

For the backward Euler method, we have

$$\left\| \frac{\partial \vec{g}}{\partial \vec{\omega}} \right\| = \left\| \frac{\partial \vec{f}}{\partial \vec{\omega}}(t_{n+1}, \vec{y}_{n+1}) \right\| \leq \frac{\lambda}{h},$$

and so (i) gives

$$h \left\| \frac{\partial \vec{f}}{\partial \vec{\omega}} \right\| \leq \lambda < 1$$

which means

$$h < \left\| \frac{\partial \vec{f}}{\partial \vec{\omega}} \right\|^{-1}.$$

So the fixed point iteration imposes a restriction on  $h$ . Can we do better?

§7.2

### Newton-Raphson method

We are trying to solve  $\vec{\omega} = h\vec{g}(\vec{\omega}) + \vec{\beta}$ . Note

$$\begin{aligned} \vec{\omega} &= h\vec{g}(\vec{\omega}^{[i]} + \vec{\omega} - \vec{\omega}^{[i]}) + \vec{\beta} \\ &= \vec{\beta} + h\vec{g}(\vec{\omega}^{[i]}) + h \frac{\partial \vec{g}}{\partial \vec{\omega}}(\vec{\omega}^{[i]}) (\vec{\omega} - \vec{\omega}^{[i]}) + \mathcal{O}(\|\vec{\omega} - \vec{\omega}^{[i]}\|^2). \end{aligned}$$

Therefore

$$\vec{\omega} - \vec{\omega}^{[i]} \approx \vec{\beta} + h\vec{g}(\vec{\omega}^{[i]}) - \vec{\omega}^{[i]} + h \frac{\partial \vec{g}}{\partial \vec{\omega}}(\vec{\omega}^{[i]}) (\vec{\omega} - \vec{\omega}^{[i]}).$$

This suggests the method

$$\underbrace{\left( \mathbf{I} - h \frac{\partial \vec{g}}{\partial \vec{\omega}}(\vec{\omega}^{[i]}) \right)}_A \underbrace{(\vec{\omega} - \vec{\omega}^{[i]})}_u \approx \underbrace{\vec{\beta} + h\vec{g}(\vec{\omega}^{[i]}) - \vec{\omega}^{[i]}}_F.$$

The Newton-Raphson method is

$$\vec{\omega}^{[i+1]} = \vec{\omega}^{[i]} - \left( \mathbf{I} - h \frac{\partial \vec{g}}{\partial \vec{\omega}}(\vec{\omega}^{[i]}) \right)^{-1} (\vec{\omega}^{[i]} - \vec{\beta} - h\vec{g}(\vec{\omega}^{[i]})).$$

Properties of Newton-Raphson:

1) For  $h > 0$  small enough, Newton-Raphson has quadratic convergence:

$$\|\vec{\omega}^{[i+1]} - \hat{\omega}\| \leq C \|\vec{\omega}^{[i]} - \hat{\omega}\|^2.$$

- 2) At every iteration, we need to compute  $\frac{\partial \vec{g}}{\partial \vec{w}}$ , which is very expensive.  
 3) At every iteration, we must solve  $AU = F$ , which is also very expensive.

Modified Newton-Raphson:

Define

$$J = \frac{\partial \vec{g}}{\partial \vec{w}}(\vec{w}^{[0]})$$

Our new scheme is

$$\vec{w}^{[i+1]} = \vec{w}^{[i]} - (I - hJ)^{-1}(\vec{w}^{[i]} - \vec{\beta} - h\vec{g}(\vec{w}^{[i]})), \quad i=0,1,2,\dots$$

This is nice because we need to compute  $J$  only once and because we can reuse  $(I - hJ)^{-1}$ . However, we lose quadratic convergence. Note that the modified Newton-Raphson method is a fixed point method:

$$\begin{aligned} \vec{w}^{[i+1]} &= \vec{w}^{[i]} - (I - hJ)^{-1}(\vec{w}^{[i]} - \vec{\beta} - h\vec{g}(\vec{w}^{[i]})) \\ \Rightarrow (I - hJ)(\vec{w}^{[i+1]} - \vec{w}^{[i]}) &= -\vec{w}^{[i]} + \vec{\beta} + h\vec{g}(\vec{w}^{[i]}) \\ \Rightarrow (I - hJ)\vec{w}^{[i+1]} - \vec{w}^{[i]} + hJ\vec{w}^{[i]} &= -\vec{w}^{[i]} + \vec{\beta} + h\vec{g}(\vec{w}^{[i]}) \\ \Rightarrow (I - hJ)\vec{w}^{[i+1]} &= h(\vec{g}(\vec{w}^{[i]}) - J\vec{w}^{[i]}) + \vec{\beta} \\ \Rightarrow \vec{w}^{[i+1]} &= h(I - hJ)^{-1}(\vec{g}(\vec{w}^{[i]}) - J\vec{w}^{[i]}) + \vec{\beta}. \end{aligned}$$

This is a fixed point method where

$$\begin{aligned} \vec{g}(\vec{w}) &= (I - hJ)^{-1}(\vec{g}(\vec{w}) - J\vec{w}), \\ \vec{\beta} &= (I - hJ)^{-1}\vec{\beta}. \end{aligned}$$

Applying theorem 7.1, we have convergence.

Recall the assumptions of theorem 7.1. If  $g$  is Lipschitz then (i) is satisfied:

$$\begin{aligned} \|\vec{g}(\vec{v}) - \vec{g}(\vec{u})\| &\leq L\|\vec{v} - \vec{u}\| \\ \frac{1}{h}\|\vec{u} - \vec{v}\| &= L\|\vec{u} - \vec{v}\| \end{aligned}$$

As  $h \rightarrow 0$  we need  $\|\vec{u} - \vec{v}\| \rightarrow 0, \rho \rightarrow 0$ .