

# Num. Stability

## Motivation:

For IVP:  $y'(t) = -10 \cdot y(t)$ ,  $0 < t < 1$ ,  $y(0) = 1$ . its sol. decays exponentially. But when we apply implicit/explicit Euler method. The num. sol. of implicit one looks good while the explicit one begins to oscillate and even explode when step length  $h \uparrow$ .

## (1) Model problem:

Fix  $\lambda \in \mathbb{C}$ .  $y'(t) = \lambda y(t)$ ,  $t_0 \leq t \leq T$ .  $y(t_0) = y_0$ .

which has sol.  $y(t) = y_0 e^{\lambda(t-t_0)}$ . (\*)

Note  $|y(t)| = |y_0| e^{\operatorname{Re}(\lambda)(t-t_0)}$ . the limit behavior depend on  $\operatorname{Re}(\lambda) > . = . < 0$  when  $t \rightarrow \infty$ .

Req: Num. sol. should imitate its exp.

decay behavior if  $\operatorname{Re}(\lambda) < 0$ .

Def: A RK method is called absolutely

stable for a  $\lambda h_k$  if it produces bad approxi.  $\sup_n |y_n| < \infty$  when applied on (\*)  
for  $\operatorname{Re}(\lambda) < 0$ . ( $h_k$  is  $k^{\text{th}}$  step length)

e.g., (Explicit Euler)

$$y_n = y_{n-1} + h f(t_{n-1}, y_{n-1}) = (1 + h\lambda) y_{n-1} \\ = \dots = (1 + h\lambda)^n y_0.$$

We call  $g(h\lambda) = 1 + h\lambda$  amplification.

factor ( $|g(h\lambda)| < 1 \Leftrightarrow |y_n| < |y_{n-1}|$ )

$\therefore$  we require:  $|g(h\lambda)| \leq 1$  for stability

e.g., (For  $\lambda = -10$ )

•  $h = \frac{1}{20}$ . it's good.

•  $h = \frac{1}{6}$ . oscillates but bad.

•  $h = \frac{1}{4}$ . explodes.

Remark: Absolutely stability only returns  
it's bad but not good approxi.

Amplification factor examples:

i) Explicit Euler  $g(\lambda h) = 1 + \lambda h$ .

ii)  $Ken$  :  $g(\lambda h) = \sum_{k=0}^2 (\lambda h)^k / k!$

iii)  $RK4$  :  $g(\lambda h) = \sum_{k=0}^4 (\lambda h)^k / k!$

$k_{max}$  : It's polynomial from Taylor expansion.

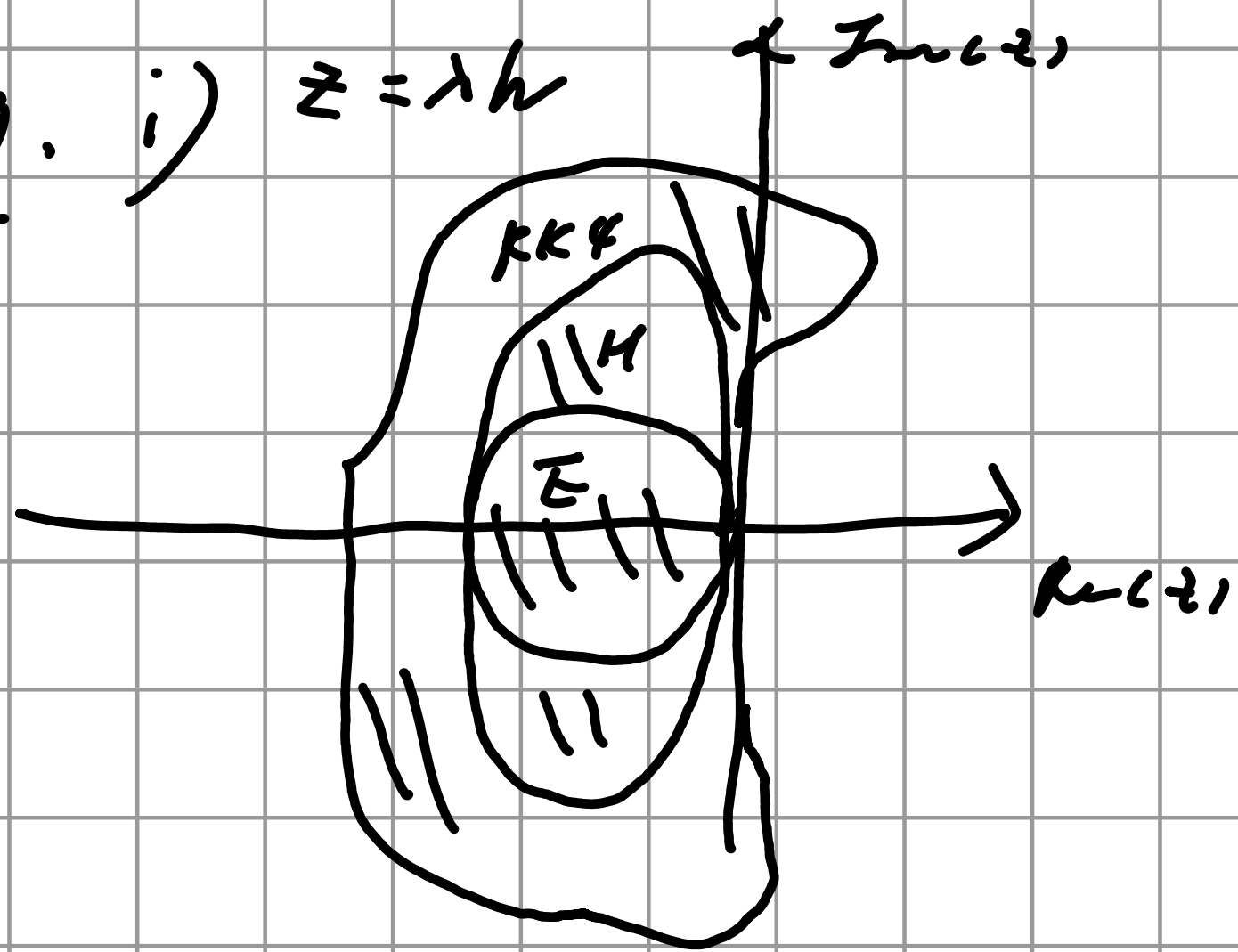
And  $|g(\lambda h)| \uparrow$  when  $|\lambda h| \uparrow$ .

Def : The stability region of a RK scheme is given by  $SR = \{\lambda h \in \mathbb{C} \mid |g(\lambda h)| \leq 1\}$ .

For  $\lambda \in \mathbb{R}$ , the stability interval is :

$RI = \{\lambda h \in \mathbb{R} \mid |g(\lambda h)| \leq 1\}$  where  $g(\lambda h)$  is amplification factor from applying scheme on  $(*)$ .

eg. i)  $z = \lambda h$



a) Explicit Euler :

$$RI = [-2, 0]$$

b)  $Ken$  :  $RI = [-2, 0]$

c)  $RK4$  :  $RI = [-2.78, 0]$ .

ii) (Nonlinear case)

For IVP  $y'(t) = \mu(y(t) - \cos(t)) - \sin(t)$ ,

$y(0) = 1$ . Note  $df/dy = \mu$ . So  $\mu$  will play the role of  $\lambda$ .

$\Rightarrow$  We need:  $\mu \leq 0$ .  $\mu h \leq 2.78$  for RK4.

Def: A method is called A-stable if  
 $L \in \mathbb{C} \mid \operatorname{Re}(L) \leq 0 \Rightarrow R(L) \leq 0$ .

Remark: Explicit scheme can't be A-stable  
 $\Rightarrow$  It motivates us to use the implicit method.

e.g. i) Implicit Euler:  $g(\lambda h) = \frac{1}{1-\lambda h}$ .



ii) Implicit Trap:  $g(\lambda h) = \frac{1+\frac{1}{2}\lambda h}{1-\frac{1}{2}\lambda h}$ .



Remark: To choose step length. We need to

consider: i) Stability (restriction on  $h$ )

ii) Accuracy:  $h$  small, accurate  $\uparrow$

iii) Efficiency:  $h$  large, efficient  $\uparrow$ .

(2) Stability func. for RK:

Apply RK method on model problem:  $\dot{x} = \lambda x$

$$k_i = \lambda \left( y_n + h \sum_j a_{ij} k_j \right), \quad y_{n+1} = y_n + h \sum_k b_k k_k$$

Let  $\tilde{k}_i = k_i / \lambda$ . We rewrite it in:

$$\underline{\tilde{k}} = y_n \cdot \begin{pmatrix} \vdots \\ 1 \end{pmatrix} + z A \cdot \underline{\tilde{k}}, \quad y_{n+1} = y_n + z \underline{b}^T \cdot \underline{\tilde{k}}$$

Where we denote  $z = h\lambda$ .

$$\Rightarrow \underline{\tilde{k}} (I_r - zA) = y_n \mathbb{I} \Rightarrow \text{solve } \underline{\tilde{k}} = (I_r - zA)^{-1} y_n \mathbb{I}.$$

$$\text{So } y_{n+1} = (1 + z \underline{b}^T (I_r - zA)^{-1} \mathbb{I}) y_n = g(z) y_n$$

$g(z)$  is amplification factor we need.

1) Explicit RK:

$$A = \begin{pmatrix} 0 & 0 \\ * & 0 \end{pmatrix} \Rightarrow A^r = 0. \quad (\text{i.e. } r\text{-nilpotent}).$$

$$\text{So } (I_r - zA)^{-1} = I + zA + \dots + (zA)^{r-1}$$

i.e.  $g(z)$  will be a polynomial of degree up to  $r$ . (So  $|g(z)| \gg 1$  if  $|z| \gg 1$ )

2) Implicit RK:

$$g(z) \text{ has form } g(z) = q(z) \cdot p(z), \quad \text{etc.}$$

$p, q$  are both polynomials of degree  $\leq r$ .

Def: One method called  $L$ -stable if it's

$A$ -stable and  $\lim_{|z| \rightarrow \infty} |g(z)| = 0$ .

Rank: i) L-stable method can have better damping property (i.e.

less oscillation when  $|z|$  large)

ii) Implicit Euler:  $|g(z)| = \frac{1}{|1-z|}$  is L-stable.

Implicit Trape.:  $|g(z)| = \left| \frac{1+\frac{1}{2}z}{1-\frac{1}{2}z} \right|$  is A-stable but not L-stable.

### 3) Application of Stab. Analy.:

Def. For IVP.  $y' = f(t, y)$ ,  $y(t_0) = y_0$ ,  $t \geq t_0$ .

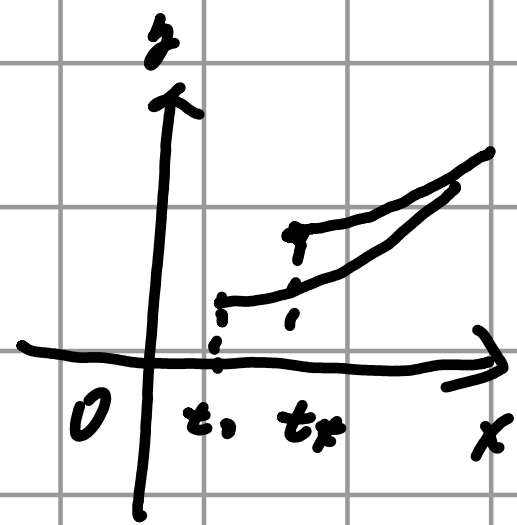
i) Its global sol. is asymptotically stab. if any sol.  $v$  of perturbed problem

$$v'(t) = f(t, v(t)), \quad v(t_x) = y(t_x) + w_x.$$

$t \geq t_x$  for some  $t_x \geq t_0$  is also global

where perturbation  $\|w_x\| \leq \delta$  &

$\|(v-y)(t)\| \rightarrow 0$  as  $t \rightarrow \infty$  holds.



ii) Let this IVP be solved by Lip. conti.

one-step method  $y_n = y_{n-1} + h_n F(h_n, t_n, y_{n-1}, y_n)$

$y_0 = y_0$ . Zts s.l.  $(y_n)$  is called numerically stable if any discrete s.l.

$(V_n)_{n \geq n^*}$  of  $V_n = V_{n-1} + h_n F(h_n, t_n, V_{n-1}, V_n)$

$V_{n^*} = y_{n^*} + W_{n^*}$ .  $n \geq n^*$ . for some  $t_{n^*} \geq t_0$

where perturbation  $\|W_{n^*}\| \leq \delta$  satisfies:

$$\|V_n - y_n\| \rightarrow 0 \text{ as } n \rightarrow \infty$$

Hypo: All eigenvalues  $\lambda(t)$  of  $f_x(t, y(t))$  are assumed to satisfy:  $\operatorname{Re}(\lambda(t)) \leq 0$ .

Prove: i) The method is asymptotically stable  
ii) The method with  $SR \subset \mathbb{C}$  is numerically stable if the step length  $h_n$  is chosen st.  $h_n \lambda(t_n) \in SR$ .  $\forall n \geq 0$ .

Remark: i)  $\lambda(t)$  will take role of  $\lambda$  in the model problem. (e.g. for  $\lambda = 1 \Rightarrow \lambda h \leq \min f_y(t, y(t))$ )

ii) It will be wrong if  $f_x$  is not diagonalizable. (counterexample exists)  
i.e. no complex system of eigenvalues

proof:

i) Gnti. problem:

$$\text{Write } w = v - y, \text{ so } w'(t) = f(t, v(t)) - f(t, y(t))$$

$$= \int_1' \frac{\partial f}{\partial s}(t, y(t) + s w(t)) ds$$

$$= \int_1' f_x(t, y(t) + s w(t)) w(t) ds$$

$$\stackrel{\text{Taylor}}{=} f_x(t, y(t)) w(t) + O(\|w(t)\|^2).$$

1) Simplification : linearization

$$\text{We get } w'(t) = f_x(t, y(t)) w(t).$$

$$w(t_*) = w_* \text{ for } t \geq t_*.$$

2) Simplification : localization.

Freeze  $t = t_*$  locally on  $y(t)$ :

$$w'(t) = f_x(t_*, y(t_*)) w(t), \quad w(t_*) = w_*, \quad t \geq t_*.$$

3) Simplification : diagonalization.

$$\exists Q, S \in \mathbb{C}. A = f_x(t, y(t_*)) \text{ have } Q A Q^{-1} = D \\ = \text{diag}(\lambda_i), \lambda_i \in \mathbb{C}.$$

$$\text{Let } \tilde{w} = Q w. \Rightarrow \tilde{w}'(t) = D \tilde{w}(t).$$

Which is reduced to modal problem.



$\operatorname{Re} \lambda_i \leq 0 \Rightarrow \widetilde{w}_i$  decays exponentially.

For  $Q$  regular ( $Q$  is invertible). We have:

$$\begin{aligned} \|W(t)\| &= \|Q^{-1} \widetilde{W}(t)\| \leq C \|\widetilde{W}(t)\| \leq C e^{\lambda(t-t_0)} \|\widetilde{W}_k\| \\ &\leq C e^{\lambda t} \|W_k\|. \quad \lambda = \max_i \operatorname{Re}(\lambda_i). \end{aligned}$$

ii) Discrete problem:

For  $g$  is amplification func. of RK method in (2) and recall it has form:

$$y_n = g(hA) y_{n-1}. \quad \text{Set } \widetilde{y}_k = Q y_k. \quad \text{We have:}$$

$$\widetilde{y}_n = Q g(hA) Q^{-1} \widetilde{y}_{n-1} = g(hD) \widetilde{y}_{n-1} \quad \text{since } g \text{ is rational func. (i.e. } p(x)/q(x) \text{).}$$

Algorithm: To solve  $y' = f(t, y)$ :

a) Evaluate  $f$  at  $(t_n, y(t_n))$ .

b) Compute eigenvalues  $\lambda_i$ .

c) Choose  $h_n$  s.t.  $h_n \lambda_i \in SR$ , for  $\operatorname{Re}(\lambda_i) \leq 0$ .

(4) Stiff problem:

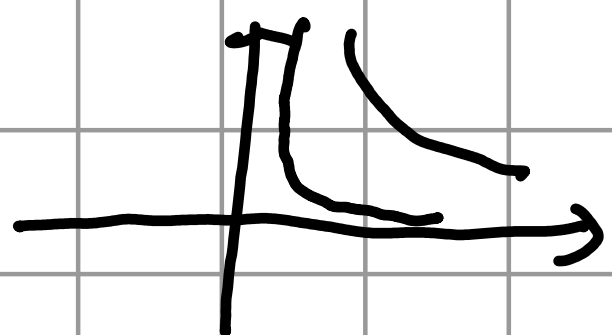
Consider the following IVP:

$$\begin{pmatrix} y_1'(t) \\ y_2'(t) \end{pmatrix} = \begin{pmatrix} -5 & 4 \\ 4 & -5 \end{pmatrix} \begin{pmatrix} y_1(t) \\ y_2(t) \end{pmatrix} \quad \begin{pmatrix} y_1(0) \\ y_2(0) \end{pmatrix} = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$$

$$\Rightarrow y_1(t) = e^{-t} + e^{-11t}, \quad y_2(t) = e^{-t} - e^{-11t}$$

But  $e^{-11t}$  decays much faster than  $e^{-t}$ .

So the behavior of s.l. will be dominated by  $e^{-t}$ .



But for num. stab.: We need to take  $e^{-11t}$  into account.

Def: i) Stiff problem is about components with different time scales

Remark: Alternatively, a IVP is stiff if the step size needed to maintain abs. stab. of expl. Euler method is much smaller than the step size needed for accuracy.

ii) The IVP is stiff along s.l. trajectory if  $\exists$  eigenvalues  $\lambda(t)$  of  $f_x(t, y(t))$

$$\text{st. } K(t) = \max_{\operatorname{Re}(\lambda(t)) < 0} |R(\lambda(t))| / \min_{\operatorname{Re}(\lambda(t)) < 0} |R(\lambda(t))| \gg 1$$

Key: i) We only consider  $\operatorname{Re}(\lambda(t)) < 0$ .

ii) There's no exact math. def.

iii) It's characterized by having components that decay in totally different speed. And we want to solve both components over long time when step length  $h$  isn't so small.

iv) Impl. methods are often A-stable so  $h$  can be chosen large to be stable. ( $\Rightarrow$  no stiff)

Key: When dealing stiff problem with some not small transient (initial error  $e_0$ )

We'd like to use the method that is L-stable (e.g. For trapezoidal: its amplification factor is  $1 + \frac{h\lambda}{2} / 1 - \frac{h\lambda}{2} \approx 1$  if  $h$  is small. So the transient will not

being fast. While implicit Euler has  $\gamma_{crit}$  is small so the transient being fast then it performs better.)

(5) Example of Impl. RK:

① Gauss method:

- i) It's based on Gaussian quadrature.
- ii) It's A-stable.
- iii) It has order  $2s$  for  $s$ -stage method.

② Radau method:

- i) It's L-stable.
- ii) It has order  $2s-1$  for  $s$ -stage method

Remark: For  $s=1$ . ①. ② both implicit Euler.