

Yaiza Pérez Tejera

# *Teoría de las estrellas de orden*

Order star theory

Trabajo Fin de Grado  
Grado en Matemáticas  
La Laguna, Julio de 2023

DIRIGIDO POR

*María Soledad Pérez Rodríguez*

*María Soledad Pérez Rodríguez*

*Análisis Matemático*

*Universidad de La Laguna*

*38200 La Laguna, Tenerife*

---

## Agradecimientos

En primer lugar agradecerle a mi tutora María Soledad Pérez por su dedicación, por todo lo que me ha enseñado, ayudado y guiado en este TFG. También agradecerle a mi familia que han estado a lo largo de toda mi carrera apoyándome en todo momento y animándome a seguir adelante. Por último pero no menos importante, a mis amigos que me han tendido una mano cuando lo he necesitado y me han hecho pasar momentos inolvidables en mi etapa universitaria. Gracias a todos, sin ustedes no hubiera sido posible.

Yaiza Pérez Tejera  
La Laguna, 10 de julio de 2023



---

## Resumen · Abstract

### *Resumen*

---

*En el presente Trabajo de Fin de Grado se ha estudiado la A-estabilidad de los métodos Runge-Kutta y, en particular, el estudio de la teoría de las estrellas de orden, que determina condiciones necesarias y suficientes para que la función de estabilidad de un método sea A-estable.*

*Para ello, hemos dividido este trabajo en tres capítulos. En el primero introducimos los métodos Runge-Kutta y algunas de sus propiedades. En el segundo capítulo estudiamos su estabilidad y definimos la función de estabilidad y la A-estabilidad. En último lugar, desarrollamos la teoría de las estrellas de orden, con la que garantizamos la A-estabilidad de los aproximantes de Padé bajo ciertas condiciones.*

**Palabras clave:** *Métodos Runge-Kutta { Estabilidad { A-estabilidad { Aproximantes de Padé { Estrellas de orden.*

### *Abstract*

---

*In this Bachelor's thesis, the A-stability of Runge-Kutta methods has been studied, with a particular focus on the theory of order stars, which determines necessary and sufficient conditions for the stability function of a method to be A-stable.*

*To do this, we have divided this work into three chapters. In the first chapter, we introduce the Runge-Kutta methods and some of their properties. In the second chapter, we study their stability and define the stability function and A-stability. Finally, we develop the theory of order stars, which guarantees the A-stability of Padé approximants under certain conditions.*

**Keywords:** *Runge-Kutta methods { Stability { A-stability { Padé approximants { Order stars.*



---

# Índice general

<b>Agradecimientos</b> .....	III
<b>Resumen/Abstract</b> .....	V
<b>Introducción</b> .....	IX
<b>1. Introducción a los métodos Runge-Kutta</b> .....	1
1.1. Introducción .....	1
1.2. Formulación de los métodos de Runge-Kutta .....	2
1.3. Método de Euler .....	4
1.4. Ejemplos de métodos RK .....	5
1.5. Existencia de la solución numérica de los RKI .....	6
1.6. Orden de los métodos Runge-Kutta .....	8
<b>2. Estabilidad de los métodos RK</b> .....	11
2.1. Ecuación test de Dahlquist .....	11
2.2. Función de estabilidad .....	12
2.3. A-estabilidad .....	15
2.4. I-estabilidad .....	16
2.5. L-estabilidad .....	18
2.6. Ejemplos .....	19
<b>3. Estrellas de orden (Order Stars)</b> .....	21
3.1. Aproximación de Padé a la función exponencial .....	21
3.2. Estrellas de orden (Order Stars) .....	23
3.3. Order Star cuando $z \neq 0$ .....	24
3.4. Order Star cuando $z \neq 1$ .....	30
3.5. Ceros y polos de $R(z)$ .....	35
3.6. Orden y estabilidad para $R(z)$ .....	41
3.7. Estabilidad de las aproximaciones de Padé .....	44
3.8. Ejemplos de aproximantes de Padé .....	45

<b>Bibliografia</b> .....	49
<b>Poster</b> .....	51

---

## Introducción

En el campo del Análisis Numérico, uno de los temas fundamentales es el estudio de los métodos de resolución de ecuaciones diferenciales que tienen aplicaciones en diversas áreas, como la física, la ingeniería, la biología, entre otras. En particular, los métodos de Runge-Kutta (RK) han demostrado ser una herramienta muy importante para aproximar soluciones numéricas con alta precisión.

Este Trabajo de Fin de grado tiene como principal objetivo estudiar las estrellas de orden, una herramienta que permite clasificar y comparar la A-estabilidad de los métodos.

En esta memoria comenzamos examinando los métodos RK y una serie de propiedades que los caracterizan. En el primer capítulo, se introducen los métodos RK, tanto explícitos como implícitos, y su orden de convergencia, viendo ejemplos de diferentes órdenes. Para ello principalmente utilizaremos [4], [5] y [6].

En el segundo capítulo se estudia la estabilidad de los métodos numéricos siguiendo [4] y [5]. La estabilidad de un método numérico tiene que ver con la manera en que los errores numéricos se propagan. Introduciremos la función de estabilidad, así como ciertas propiedades que satisface ésta demostrando que es una función racional que aproxima a la función exponencial. Se define entonces el concepto de la A-estabilidad de los métodos o, lo que es lo mismo, que la función de estabilidad está acotada por 1 en el semiplano complejo negativo, lo que garantiza que el método será estable sobre problemas lineales de coeficientes constantes. También se definen otros conceptos relacionados como la I-estabilidad y la L-estabilidad, viendo ejemplos de métodos que satisfacen estas propiedades.

En el último capítulo, abordamos el estudio de la teoría de las estrellas de orden (*order stars*), donde principalmente desarrollamos parte del Capítulo

IV.4 de [5]. Esta teoría se estudia en general para profundizar en la relación entre los aproximantes racionales a la exponencial y dicha exponencial. Con ello, por una parte, si tenemos una aproximación racional a  $e^z$  podremos deducir si es A-estable o no y, por otro lado, ayuda a construir aproximaciones racionales que sean A-estables. En particular, se estudian los aproximantes de Padé a  $e^z$ , que alcanzan el orden de aproximación mayor posible (alrededor de  $z = 0$ ), entre todas las aproximaciones racionales a la exponencial con los grados del numerador y del denominador prefijados.

El objetivo de la teoría desarrollada permitirá demostrar el principal teorema de este trabajo que garantiza la A-estabilidad de los aproximantes de Padé a la exponencial.

## Introducción a los métodos Runge-Kutta

En este primer capítulo introduciremos los métodos Runge-Kutta, una familia de algoritmos que se utilizan ampliamente para la solución numérica de EDOs. Haremos hincapié en algunas de sus propiedades y visualizaremos algunos ejemplos. Por último, se demostrará el teorema que garantiza la definición de los métodos RK implícitos.

Las principales propiedades sobre ecuaciones diferenciales las hemos extraído de [6] y [1]. Para introducir los métodos RK y sus propiedades hemos seguido [4] y [5].

### 1.1. Introducción

Consideremos el problema de valor inicial (PVI)

$$y'(t) = f(t, y(t)), \quad y(t_0) = y_0, \quad t \in [t_0, t_f] \quad (1.1)$$

con  $f : [t_0, t_f] \times \Omega \rightarrow \mathbb{R}^m$ , donde  $\Omega \subset \mathbb{R}^m$  es un abierto no vacío.

Una **solución** de (1.1) es una función continuamente diferenciable en algún intervalo  $I$  que contiene a  $t_0$  y que verifica (1.1).

**Teorema 1** *Teorema de existencia*[6]

Sea  $f$  continua en  $[t_0, t_f] \times \Omega$ . Para cada  $y_0 \in \Omega$  existe una solución de (1.1).

**Definición 1** Una función  $f : [t_0, t_f] \times \Omega \rightarrow \mathbb{R}^m$  se dice **Lipschitz** en  $[t_0, t_f] \times \Omega$  si existe una constante  $L > 0$  tal que

$$\|f(s, u) - f(s, v)\| \leq L\|u - v\|, \quad \forall s \in [t_0, t_f], \quad \forall u, v \in \Omega.$$

**Definición 2** Una función  $f : [t_0, t_f] \times \Omega \rightarrow \mathbb{R}^m$  se dice **localmente Lipschitz** si para cada  $(t, y) \in [t_0, t_f] \times \Omega$  existen entornos abiertos  $I \subset [t_0, t_f], U \subset \Omega$  y una constante  $L > 0$  tal que  $(t, y) \in I \times U$  y

$$\|f(s, u) - f(s, v)\| \leq L \|u - v\|, \quad \forall s \in I, \forall u, v \in U \quad (1.2)$$

**Observación 1** Si  $f$  es una función Lipschitz entonces  $f$  es localmente Lipschitz.

**Teorema 2 [6]** Sea  $f : [t_0, t_f] \times \Omega \rightarrow \mathbb{R}^m$  continua y localmente Lipschitz con respecto a su segunda componente. Entonces, para cada  $y_0 \in \Omega$ , existe una única solución de (1.1).

En este trabajo estamos interesados en la aproximación numérica de las soluciones de PVI's (1.1). En particular, la aproximación dada por los métodos Runge-Kutta.

Para simplificar la explicación, en lo que sigue consideraremos que la función derivada  $f$  está definida, es continua y Lipschitz respecto de  $y$  en  $I \times \mathbb{R}^m$  donde  $I \subset \mathbb{R}$  es cualquier intervalo cerrado.

En estos casos utilizaremos como normas en  $\mathbb{R}^m$ ,

$$\|v\|_1 = \max_{i=1, \dots, m} |v_i| \quad \text{ó} \quad \|v\|_2 = \sqrt{\sum_{i=1}^m v_i^2}$$

donde son la norma uniforme y la norma  $l_2$  respectivamente.

## 1.2. Formulación de los métodos de Runge-Kutta

Dado el PVI (1.1), para resolverlo numéricamente se toma una partición del intervalo de integración  $[t_0, t_f]$

$$P_N = \{t_0 < t_1 < \dots < t_N = t_f\} \quad (1.3)$$

y se calculan una serie de aproximantes  $y_0, y_1, \dots, y_N$  tal que

$$y_n = y(t_n), \quad n = 0, 1, \dots, N$$

mediante el siguiente algoritmo.

**Definición 3** Dado un entero  $s \geq 1$  (numero de etapas) el método que, dada una aproximación  $y_n$  a la solución del PVI (1.1) en  $t_n$  nos da una aproximación  $y_{n+1}$  a la solución en  $t_{n+1} = t_n + h$  mediante la fórmula:

$$k_i = f(t_n + c_i h, y_n + h \sum_{j=1}^s a_{ij} k_j), \quad i = 1, \dots, s \tag{1.4}$$

$$y_{n+1} = y_n + h \sum_{i=1}^s b_i k_i \tag{1.5}$$

se llama método **Runge-Kutta de  $s$  etapas (RK)**.

Para manejar los coeficientes que definen el método, se define la **tabla de Butcher** asociada al RK (1.4) como:

$$\begin{array}{c|cccc} c_1 & a_{11} & a_{12} & \dots & a_{1s} \\ c_2 & a_{21} & a_{22} & \dots & a_{2s} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ c_s & a_{s1} & a_{s2} & \dots & a_{ss} \\ \hline & b_1 & b_2 & \dots & b_s \end{array}$$

donde la matriz  $A = (a_{ij})_{i,j=1}^s$  se llama **matriz de coeficientes**, el vector  $c = (c_1, c_2, \dots, c_s)^T$  es el vector de **nodos** y  $b = (b_1, b_2, \dots, b_s)^T$  es el vector de **pesos**.

Como es usual, a partir de ahora supondremos que los pesos y los nodos satisfacen:

$$\sum_{i=1}^s b_i = 1, \quad c_i = \sum_{j=1}^s a_{ij}, \quad i = 1, \dots, s. \tag{1.6}$$

Estas expresiones también pueden escribirse de forma vectorial como  $b^T e = 1$  y  $Ae = c$  donde  $e = (1, 1, \dots, 1)^T$ , por lo que denotaremos como  $RK(A, b)$  al método RK con matriz  $A$ , pesos  $b$  y nodos  $c = Ae$ .

Según la forma de la matriz  $A$  los métodos RK se pueden clasificar en dos grupos:

1. Cuando la matriz  $A$  es triangular inferior estricta, el método se dice **explícito (RKE)**, pues sus etapas  $k_i, i=1, \dots, s$  se obtienen de forma recursiva, es decir,

$$\begin{cases} k_i = f(t_n + c_i h, y_n + h \sum_{j=1}^{i-1} k_j), & 1 \leq i \leq s \\ y_{n+1} = y_n + h \sum_{i=1}^s b_i k_i \end{cases} \tag{1.7}$$

2. Cuando  $a_{ij} \neq 0$  para algún  $j > i$ , el método se llama **implícito (RKI)**. En estos métodos necesitaremos resolver el sistema implícito (1.4) de dimensión  $s \times m$  para calcular sus etapas  $f_k, g_{i=1}^s$ .

### 1.3. Método de Euler

Euler fue el primero que propuso aproximar la solución de un PVI. El método de Euler también es el más simple de los métodos numéricos para resolver un PVI's (1.1).

Consideramos la partición  $P_N$  dada en (1.3) y denotamos  $h_j = t_{j+1} - t_j$ ,  $j = 0, \dots, N-1$ . Por la definición de derivada para cada,  $k = 0, 1, \dots, N-1$ ,

$$y'(t_k) = \lim_{h \rightarrow 0} \frac{y(t_k + h) - y(t_k)}{h},$$

por lo que esperamos que para  $h_k$  "suficientemente pequeño"

$$y'(t_k) \approx f(t_k, y(t_k)) \approx \frac{y(t_k + h_k) - y(t_k)}{h_k}.$$

Por tanto, como  $t_k + h_k = t_{k+1}$ :

$$y(t_{k+1}) \approx y(t_k) + h_k f(t_k, y(t_k)), \quad k = 0, 1, \dots, N-1.$$

Como estamos buscando unos aproximantes  $y_k \approx y(t_k)$ , esto nos sugiere el cálculo de los  $y_k$  mediante la recurrencia,

$$y_0 = y(t_0), \quad y_{k+1} = y_k + h_k f(t_k, y_k), \quad k = 0, 1, \dots, N-1, \quad (1.8)$$

Esta fórmula (1.8) se conoce como **método de Euler**, y es un caso particular de RK (1.4)-(1.5) con  $s = 1$ ,  $A = 0$  y  $b = 1$ .

#### **Teorema 3** Teorema de convergencia del método de Euler:[1]

Consideramos el PVI (1.1) con  $f$  continua y Lipschitz respecto de  $y$ . Si  $y_n, g_{n=0}^N$  es la solución dada por el método de Euler (1.8) sobre la partición (1.3), y la solución exacta  $y(t)$  verifica

$$\|y''(t)\| \leq Y_2, \quad \forall t \in [t_0, t_f]$$

entonces

$$\max_{t_n \in P} \|y_n - y(t_n)\| \leq \left( \frac{Y_2 (e^{L(t_f - t_0)} - 1)}{2L} \right) h_{max}$$

donde  $h_{max} = \max h_j$  y  $L$  es la constante de Lipschitz de  $f$ .

### 1.4. Ejemplos de métodos RK

A partir del método de Euler y con el desarrollo de las fórmulas de cuadratura, a finales del siglo *XIX* varios autores comenzaron a proponer otros métodos de un paso para resolver PVI, que más tarde se incluyeron en la clase de métodos RK. A continuación, pondremos algunos ejemplos:

1. Método de Runge (1905): (4 etapas)

$$\left\{ \begin{array}{l} k_1 = f_n \\ k_2 = f\left(t_n + \frac{h}{2}, y_n + \frac{h}{2}k_1\right) \\ k_3 = f(t_n + h, y_n + hk_2) \\ k_4 = f(t_n + h, y_n + hk_3) \\ y_{n+1} = y_n + h \left( \frac{k_1}{6} + \frac{2}{3}k_2 + \frac{1}{6}k_4 \right) \end{array} \right. \quad \begin{array}{c|cccc} 0 & 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ \hline & 1/6 & 2/3 & 0 & 1/6 \end{array} \quad (1.9)$$

2. Método de Kutta (1905):

$$\left\{ \begin{array}{l} k_1 = f_n \\ k_2 = f\left(t_n + \frac{h}{2}, y_n + \frac{h}{2}k_1\right) \\ k_3 = f\left(t_n + \frac{h}{2}, y_n + \frac{h}{2}k_2\right) \\ k_4 = f(t_n + h, y_n + hk_3) \\ y_{n+1} = y_n + h \left( \frac{k_1}{6} + \frac{k_2}{3} + \frac{k_3}{3} + \frac{k_4}{6} \right) \end{array} \right. \quad \begin{array}{c|cccc} 0 & 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 & 0 \\ 1/2 & 0 & 1/2 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ \hline & 1/6 & 1/3 & 1/3 & 1/6 \end{array} \quad (1.10)$$

3. Método de Euler implícito

$$y_1 = y_0 + hf(x_1, y_1) \quad \frac{1|1}{1} \quad (1.11)$$

4. Regla del punto medio implícito

$$\begin{array}{l} k_1 = f\left(x_0 + \frac{h}{2}, y_0 + \frac{h}{2}k_1\right) \\ y_1 = y_0 + hk_1 \end{array} \quad \frac{1/2|1/2}{1} \quad (1.12)$$

5. Regla trapezoidal

$$\begin{array}{l}
 k_1 = f(t_n, y_n) \\
 k_2 = f(t_n + h, y_n + \frac{h}{2}(k_1 + k_2)) \\
 y_{n+1} = y_n + \frac{h}{2}(k_1 + k_2)
 \end{array}
 \quad
 \begin{array}{l}
 0 \\
 1 \mid 1/2 \ 1/2 \\
 \hline
 1/2 \ 1/2
 \end{array}
 \quad (1.13)$$

Como se puede observar los dos primeros son explícitos mientras que los demás son implícitos.

### 1.5. Existencia de la solución numérica de los RKI

Para calcular la aproximación  $y_{n+1}$  en cada método implícito (1.4)-(1.5), necesitamos asegurar que el sistema implícito (1.4) de las etapas tiene solución única.

**Teorema 4 [4]** Sea  $f : [t_0, t_f] \times \mathbb{R}^m \rightarrow \mathbb{R}^m$  una función continua que satisface una condición de Lipschitz respecto a  $y$ , siendo  $L$  la constante de Lipschitz. Si tomamos

$$\tilde{h} < \frac{1}{L \max_i \sum_j |a_{ij}|}, \quad (1.14)$$

para todo  $h \in [0, \tilde{h}]$  existe una solución única del sistema implícito (1.4) que se puede obtener por iteración funcional. Si  $f$  es  $p$  veces continuamente diferenciable, las funciones  $k_i = k_i(h)$  también son de  $C^p([0, \tilde{h}])$ .

*Demostración:* Sea  $K = (k_1, k_2, \dots, k_s)^T \in \mathbb{R}^{sm}$  y consideramos la norma  $\|K\| = \max_{1 \leq i \leq s} \|k_i\|$  donde  $\|\cdot\|$  es la norma considerada en  $\mathbb{R}^m$ . El sistema implícito (1.4) se puede reescribir como  $K = F(K)$  donde  $F(K) = (F_1(K), \dots, F_s(K))^T$ , siendo

$$F_i(K) = f(t_n + c_i h, y_n + h \sum_{j=1}^s a_{ij} k_j), \quad 1 \leq i \leq s, \quad h \in [0, \tilde{h}].$$

Considérense  $K^1, K^2 \in \mathbb{R}^{sm}$ . Para cada  $i = 1, 2, \dots, s$ ,

$$\|F_i(K^1) - F_i(K^2)\| = \|f(t_n + c_i h, y_n + h \sum_{j=1}^s a_{ij} k_j^1) - f(t_n + c_i h, y_n + h \sum_{j=1}^s a_{ij} k_j^2)\|$$

$$\leq L h \sum_{j=1}^s |a_{ij}| \|k_j^1 - k_j^2\| \leq h L \sum_{j=1}^s |a_{ij}| \|k_j^1 - k_j^2\|$$

Por tanto,

$$\|F(K^1) - F(K^2)\| = \max_{1 \leq i \leq s} \|F_i(K^1) - F_i(K^2)\|$$

$$\begin{aligned} & \max_{1 \leq i \leq s} hL \sum_{j=1}^s ja_{ij} k_j^1 \quad k_j^2 \\ & \tilde{h}L \left[ \max_{1 \leq i \leq s} \sum_{j=1}^s ja_{ij} \right] \quad K^1 \quad K^2 \end{aligned}$$

Por hipótesis tenemos que  $\tilde{h}L \max_{1 \leq i \leq s} \sum_{j=1}^s ja_{ij} < 1$ , por lo que para todo  $h \in [0, \tilde{h}]$ ,

$$\|F(K^1) - F(K^2)\| \leq M \|K^1 - K^2\|, \quad M = \tilde{h}L \max_{1 \leq i \leq s} \sum_{j=1}^s ja_{ij} < 1.$$

En consecuencia,  $F$  es contractiva. Por el Teorema del Punto Fijo, existe solución única del sistema  $K = F(K)$  y dicha solución se puede aproximar por iteración funcional

$$K^{(\nu+1)} = F(K^{(\nu)}), \quad \nu = 0, 1, \dots, \quad \text{dado } K^{(0)}.$$

En otras palabras, tomando un valor inicial  $K^{(0)} = (k_1^{(0)}, \dots, k_s^{(0)})$  se obtiene para cada  $i = 1, \dots, s$ , la aproximación a la solución de (1.4) mediante la iteración:

$$k_i^{(\nu+1)} = f \left( t_n + c_i h, y_n + h \sum_{j=1}^s a_{ij} k_j^{(\nu)} \right), \quad \nu = 0, 1, 2, \dots$$

Para demostrar que la ecuación (1.4) define a las etapas  $k_i$  como funciones derivables de  $h$ , vamos a aplicar el Teorema de la Función Implícita [7]. Para ello reescribimos (1.4) como

$$\phi(h, K) = K - F(K) = 0$$

En consecuencia,

$$\frac{\partial \phi(h, K)}{\partial K} = I - \frac{\partial F}{\partial K}, \quad \frac{\partial F}{\partial K} = \frac{\partial (F_1(K), \dots, F_s(K))}{\partial (k_1, k_2, \dots, k_s)}$$

Para cada  $i, l = 1, 2, \dots, s$ :

$$\frac{\partial F_i(k_1, \dots, k_s)}{\partial k_l} = \frac{\partial}{\partial k_l} \left( f(t_n + c_i h, y_n + h \sum_{j=1}^s a_{ij} k_j) \right) \quad (1.15)$$

$$= \frac{\partial f}{\partial y} (t_n + c_i h, y_n + h \sum_{j=1}^s a_{ij} k_j) \quad h a_{il} \quad (1.16)$$

Entonces, si  $h = 0$ ,  $\partial F_i / \partial k_l = 0$ ,  $\delta_{il}$ ,  $l = 1, \dots, s$ , o bien  $\partial F / \partial K = 0$ , lo que implica que

$$\frac{\partial \phi}{\partial K}(0, K) = I$$

En consecuencia, por el Teorema de la Función Implícita  $\phi(h, K) = 0$  define implícitamente a  $K$  como función de  $h$ ,  $K = K(h)$ , y además, como  $f$  es  $p$  veces diferenciable con continuidad,  $K(h)$  es de clase  $C^p([0, \tilde{h}])$ .

## 1.6. Orden de los métodos Runge-Kutta

Hemos visto que los métodos RK están bien definidos. De forma natural, tenemos que ver si las aproximaciones que dan los métodos realmente aproximan a la solución del PVI (1.1).

**Definición 4** Un método  $RK(A, b)$  es de **orden  $p \geq 1$**  si para problemas (1.1) con  $f \in C^p([t_0, t_f], \Omega)$  y lipschitz, se verifica

$$\|y(t_0 + h) - y_1\| \leq Kh^{p+1}, \quad h \in 0^+, \quad (1.17)$$

es decir, si la serie de Taylor para la solución exacta  $y(t_0 + h)$  y la aproximación  $y_1$  coinciden hasta el término  $h^p$ , cuando  $h \in 0^+$ .

**Teorema 5** Un  $RK(A, b)$  tiene al menos orden 1,  $b^T e = 1$

*Demostración:* Para tener orden 1 debe verificarse que

$$\|y(t_0 + h) - y_1\| \leq Kh^2, \quad h \in 0$$

El desarrollo de Taylor cuando  $h \in 0$ , es

$$y(t_0 + h) = y(t_0) + hy'(t_0) + O(h^2) = y_0 + hf(t_0, y_0) + O(h^2) \quad (1.18)$$

Además la solución en  $t_0 + h$  es  $y_1 = y_0 + \sum_{i=1}^s b_i K_i$  donde las etapas del método son:

$$K_i = f(t_0 + c_i h, y_0 + h \sum_{j=1}^{i-1} a_{ij} K_j) = f(t_0, y_0) + O(h), \quad i = 1, 2, \dots, s.$$

Por tanto,

$$y_1 = y_0 + h \sum_{i=1}^s b_i (f(t_0, y_0) + O(h)) = y_0 + h \left( \sum_{i=1}^s b_i \right) f(t_0, y_0) + O(h^2)$$

Comparando con (1.18)

$$y(t_0 + h) - y_1 = h \left( 1 - \sum_{i=1}^s b_i \right) f(t_0, y_0) + O(h^2)$$

Luego para tener orden 1 es necesario y suficiente que

$$1 - \sum_{i=1}^s b_i = 0 \quad \text{y} \quad \sum_{i=1}^s b_i = 1.$$

Existen formas de saber si un método RK determinado va a converger y con qué orden. Un requerimiento fundamental para poder garantizar dicha convergencia es que el método propague bien los errores a medida que se van calculando los aproximantes, es decir, garantizar que el método es estable. En los siguientes capítulos nos centraremos en el estudio de la estabilidad.



---

## Estabilidad de los métodos RK

La estabilidad es un aspecto crucial al trabajar con métodos numéricos para resolver ecuaciones diferenciales. En el caso de los métodos Runge-Kutta, la estabilidad se convierte en un factor determinante para garantizar la fiabilidad de las soluciones aproximadas. Lo primero que se ve cuando se estudia dicha estabilidad es la diferencia que hay entre los métodos explícitos y los implícitos.

Como ejemplo lustrativo de dicha diferencia de estabilidad entre los métodos explícitos e implícitos comenzaremos comparando la aplicación del método de Euler explícito con el implícito sobre un problema que parece sencillo, la ecuación test de Dahlquist. Para el desarrollo de este capítulo hemos extraído principalmente la información de [5, Cap.IV.3].

### 2.1. Ecuación test de Dahlquist

Consideremos la **ecuación test de Dahlquist**

$$y' = \lambda y, \quad \operatorname{Re} \lambda < 0 \quad (2.1)$$

El método Euler explícito sobre (2.1) con paso fijo  $h = h_j, \delta_j$ , nos da las aproximaciones

$$y_n = y_{n-1} + hf(t_{n-1}, y_{n-1}) = y_{n-1} + h\lambda y_{n-1} = R(h\lambda)y_{n-1}, \quad R(z) = 1 + z$$

por lo que, iterando,

$$y_n = R(h\lambda)^n y_0$$

Por tanto, si  $h$  es tal que  $|R(h\lambda)| > 1$ , la solución numérica explotará a medida que calculemos más pasos. Así que para que el método dé una aproximación satisfactoria, necesariamente se tiene que tomar  $h$  suficientemente pequeño para que

$$|R(h\lambda)| < 1, \quad |1 + h\lambda| < 1, \quad (2.2)$$

o, lo que es lo mismo, que  $h\lambda$  caiga dentro del círculo centrado en  $-1$  de radio  $1$ . Si por ejemplo,  $\lambda = -10^3$ , esto exige una restricción de paso

$$|1 - h\lambda| \leq 1, \quad |1 + h\lambda| \leq 1, \quad h \leq 2 \cdot 10^{-3}$$

Esto supone que, si integramos el problema para  $t \in [0, 1]$ , el método tendrá que dar al menos 500 pasos para obtener una aproximación razonable.

En cambio, sobre Euler implícito (2.1)

$$y_n = y_{n-1} + hf(t_n, y_n) = y_{n-1} + h\lambda y_n$$

$$(1 - h\lambda)y_n = y_{n-1} \Rightarrow y_n = R(h\lambda)y_{n-1}, \quad R(z) = \frac{1}{1 - z}$$

En este caso,

$$|R(h\lambda)|^j = \frac{1}{|1 - h\lambda|^j} < 1, \quad |1 - h\lambda| > 1$$

Es decir,  $h\lambda$  tiene que estar fuera del círculo de centro  $1$  y radio  $1$ , lo que se da para cualquier  $h > 0$ , ya que  $Re\lambda < 0$ . Por tanto, no hay ninguna restricción de paso por estabilidad en este método.

Como ya veremos, lo que le pasa al método de Euler explícito se dará también para los métodos RKE en mayor o menor medida. Para verlo, tendremos que estudiar como son las funciones  $R(z)$  de los RK.

## 2.2. Función de estabilidad

Consideremos un método  $RK(A, b)$  de  $s$  etapas:

$$\begin{cases} k_i = f\left(t_n + c_i h, y_n + h \sum_{j=1}^s a_{ij} k_j\right), & 1 \leq i \leq s \\ y_{n+1} = y_n + h \sum_{i=1}^s b_i k_i \end{cases} \quad (2.3)$$

**Proposición 1** *El método  $RK(A, b)$  (2.3) aplicado a la ecuación de Dahlquist (2.1) resulta:*

$$y_{n+1} = R(h\lambda)y_n, \quad R(z) = 1 + zb^T(I - zA)^{-1}e \quad (2.4)$$

donde  $e = (1, \dots, 1)^T$  e  $I$  es la matriz identidad de dimensión  $s$ .

La función  $R(z)$  se llama **función estabilidad** del método (2.3).

*Demostración:* Como  $f(t, y) = \lambda y \in \mathbb{R}$ , las etapas  $k_i$  de (2.3) resultan

$$k_i = \lambda \left( y_n + h \sum_{j=1}^s a_{ij} k_j \right), \quad i = 1, \dots, s$$

Matricialmente, si  $K = (k_1, \dots, k_s)^T$

$$K = \lambda y_n e + \lambda h A K \quad (I - \lambda h A) K = \lambda y_n e$$

Luego para  $h$  suficientemente pequeño,  $I - \lambda h A$  es invertible y

$$K = \lambda (I - \lambda h A)^{-1} y_n e \quad (2.5)$$

Por tanto,

$$y_{n+1} = y_n + h \sum b_i K_i = y_n + h b^T K$$

$$y_{n+1} = y_n + h b^T (\lambda (I - h \lambda A)^{-1} y_n e) = y_n + h \lambda (b^T (I - h \lambda A)^{-1} e) y_n$$

Luego, denotando  $z = h \lambda$  nos queda:

$$y_{n+1} = (1 + z b^T (I - z A)^{-1} e) y_n = R(h \lambda) y_n$$

con lo que se demuestra (2.4).

**Proposición 2** *La función de estabilidad (2.4) satisface*

$$R(z) = \frac{\det(I - zA + z e b^T)}{\det(I - zA)} \quad (2.6)$$

*Demostración:* Si llamamos

$$g_i = y_n + h \sum_{j=1}^s a_{ij} k_j, \quad i = 1, \dots, s, \quad g = (g_1, \dots, g_s)^T,$$

usando la misma notación matricial que en la proposición anterior, se tiene que  $g = e y_n + h A K$ .

Por (2.5) tenemos,

$$\begin{aligned} g &= e y_n + h \lambda A (I - h \lambda A)^{-1} y_n e = (I + h \lambda A (I - h \lambda A)^{-1}) e y_n \\ &= ((I - h \lambda A) + h \lambda A) (I - h \lambda A)^{-1} e y_n = (I - h \lambda A)^{-1} e y_n \end{aligned}$$

Entonces,

$$(I - z A) g = e y_n \quad (2.7)$$

Por otro lado,

$$\begin{aligned} y_{n+1} &= y_n + hb^T K = y_n + h\lambda b^T (I - h\lambda A)^{-1} e y_n = \\ &= y_n + h\lambda b^T (I - h\lambda A)^{-1} e y_n = y_n + zb^T g \end{aligned}$$

Luego,

$$zb^T g + y_{n+1} = y_n \quad (2.8)$$

Agrupando (2.7) y (2.8) obtenemos un sistema lineal de  $s + 1$  ecuaciones con incógnitas  $\bar{f}g, y_{n+1}g$  de la forma:

$$\begin{pmatrix} I & zA & \mathbf{0} \\ zb^T & 1 \end{pmatrix} \begin{pmatrix} g \\ y_{n+1} \end{pmatrix} = \begin{pmatrix} e y_n \\ y_n \end{pmatrix} \quad (2.9)$$

Como la matriz de coeficientes del sistema es triangular inferior el determinante será  $\det(I - zA)$ , por lo que, resolviendo el sistema para  $y_{n+1}$ , como  $y_n \in \mathbb{R}$ ,

$$y_{n+1} = \frac{\begin{vmatrix} I & zA & e y_n \\ zb^T & y_n \end{vmatrix}}{\det(I - zA)} = \frac{\begin{vmatrix} I & zA & e \\ zb^T & 1 \end{vmatrix}}{\det(I - zA)} y_n \quad (2.10)$$

Por otro lado, si denotamos  $M = I - zA = (m_{ij})_{i,j=1}^s$ , aplicando operaciones elementales,

$$\begin{vmatrix} I & zA & e \\ zb^T & 1 \end{vmatrix} = \begin{vmatrix} m_{11} & m_{12} & \dots & m_{1s} & 1 \\ m_{21} & m_{22} & \dots & m_{2s} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ m_{s1} & m_{s2} & \dots & m_{ss} & 1 \\ zb_1 & zb_2 & \dots & zb_s & 1 \end{vmatrix} = \quad (2.11)$$

$$= \begin{vmatrix} m_{11} + zb_1 & m_{12} + zb_2 & \dots & m_{1s} + zb_s & 0 \\ m_{21} + zb_1 & m_{22} + zb_2 & \dots & m_{2s} + zb_s & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ m_{s1} + zb_1 & m_{s2} + zb_2 & \dots & m_{ss} + zb_s & 0 \\ zb_1 & zb_2 & \dots & zb_s & 1 \end{vmatrix} = \begin{vmatrix} m_{11} + zb_1 & \dots & m_{1s} + zb_s \\ \vdots & \ddots & \vdots \\ m_{s1} + zb_1 & \dots & m_{ss} + zb_s \end{vmatrix} = \det(M + zeb^T) \quad (2.12)$$

y se obtiene (2.6).

Como consecuencia, desarrollando los determinantes de orden  $s$  de (2.6) se tiene que  $R(z)$  es una función racional de  $z$ , esto es

$$R(z) = \frac{P(z)}{Q(z)} \quad (2.13)$$

cuyos grados,  $\deg P = k$  y  $\deg Q = j$  respectivamente, son ambos menores o iguales a  $s$ .

Además, esta función racional es una aproximación a la función exponencial.

En lo que sigue  $P(z)$  y  $Q(z)$  denotarán el numerador y el denominador, respectivamente, de una función racional  $R(z)$  y  $(k, j)$  denotarán sus grados respectivos.

**Proposición 3** *Si el método es de orden  $p \geq 1$  entonces,*

$$e^z - R(z) = Cz^{p+1} + O(z^{p+2}) \quad \text{cuando } z \neq 0, \quad C \neq 0. \quad (2.14)$$

Como consecuencia,  $R(z)$  es una aproximación racional de  $e^z$  de orden  $p$ .

*Demostración:* Consideramos el PVI  $y' = \lambda y, \quad y(0) = 1$ . Su solución exacta en  $t_1 = h$  es  $y(h) = 1 + \lambda h + \frac{1}{2} \lambda^2 h^2 + \dots = e^{\lambda h} = e^z$  y la solución numérica en  $t_1 = h$  dada por el método  $RK(A, b)$  es  $y_1 = R(\lambda h) y_0 = R(z)$

Si el método es de orden  $p$ :

$$y(t_1) - y_1 = \tilde{C}h^{p+1} + O(h^{p+2}), \quad h \neq 0$$

Por tanto,

$$e^z - R(z) = \tilde{C}h^{p+1} + O(h^{p+2}) = Cz^{p+1} + O(z^{p+2}), \quad z \neq 0,$$

para cierto  $C$  independiente de  $z$ .

### 2.3. A-estabilidad

**Definición 5** *Un método RK, cuyo dominio de estabilidad*

$$S = \{z \in \mathbb{C} \mid |R(z)| \leq 1\} \quad (2.15)$$

satisface

$$\mathbb{C} := \{z \mid \operatorname{Re} z \leq 0\} \cap S$$

se denomina **A-estable**.

**Proposición 4**

$$\text{un } RK(A, b) \text{ es A-estable, } \begin{cases} R(z) \text{ es analítica en } \operatorname{Re} z < 0 \\ |R(iy)| \leq 1, \forall y \in \mathbb{R} \end{cases}$$

*Demostración:* Si existe  $\tilde{z} \in \mathbb{C}$  con  $Q(\tilde{z}) = 0$ ,  $R(z)$  no está acotada. Por tanto para ser A-estable  $R(z)$  tiene que ser analítica en  $\mathbb{C}$ . Además, por el principio del máximo, si  $R(z)$  es analítica en  $\mathbb{C}$ , su máximo se alcanza en su frontera, es decir, el eje imaginario.

**Proposición 5** *Ningun metodo RK explícito puede ser A-estable.*

*Demostración:* Sabemos que la matriz  $A$  en los métodos explícitos es triangular inferior estricta, por lo tanto  $\det(I - zA) = 1$ , entonces, por (2.6) necesariamente  $R(z)$  es un polinomio de grado  $s$ . Por eso cuando  $z \rightarrow \infty$  no puede estar acotado, así que no puede ser A-estable.

## 2.4. I-estabilidad

La condición  $|R(iy)| \leq 1$ ,  $y \in \mathbb{R}$ , vista en la Proposición (4), suele llamarse **I-estabilidad**.

Estudiando mejor esta condición, tenemos

$$|R(iy)|^2 \leq 1 \iff \frac{|P(iy)|^2}{|Q(iy)|^2} \leq 1 \iff |P(iy)|^2 - |Q(iy)|^2 \leq 0 \iff |Q(iy)|^2 - |P(iy)|^2 \geq 0$$

Por tanto, si denotamos

$$E(y) = |Q(iy)|^2 - |P(iy)|^2$$

la I-estabilidad es equivalente a que

$$E(y) \geq 0, \quad y \in \mathbb{R}. \quad (2.16)$$

**Proposición 6**  *$E(y)$  es un polinomio de potencias pares con*

$$\deg E = 2 \max(\deg P, \deg Q).$$

*Si  $R(z)$  es una aproximación del orden de  $p$  a la exponencial, entonces*

$$E(y) = O(y^{p+1}), \quad y \rightarrow 0$$

*Demostración:*

$$E(y) = Q(iy)Q(-iy) - P(iy)P(-iy) = E(-y), \quad y \in \mathbb{R},$$

por lo que  $E(y)$  es un polinomio par, es decir, solo tiene potencias pares. Por otro lado, por (2.14) para  $z = iy$  y usando que  $je^{iy}j = 1$ ,

$$je^{iy} \frac{jP(iy)}{jQ(iy)} = 1 \quad \frac{jP(iy)}{jQ(iy)} = O(y^{p+1})$$

de donde se deduce que

$$jQ(iy) - jP(iy) = O(y^{p+1})$$

En consecuencia,

$$E(y) = jQ(iy)^2 - jP(iy)^2 = (jQ(iy) + jP(iy))(jQ(iy) - jP(iy)) = O(y^{p+1})$$

ya que  $(jQ(iy) + jP(iy))$  está acotado.

La siguiente propiedad, que aplicaremos en el siguiente capítulo, se deduce de lo anterior.

**Proposición 7** Una función racional  $R(z)$  (2.13) de orden  $p = 2j - 2$  es I-estable si y solo si  $jR(1) = 1$ .

*Demostración:* Por la Proposición 6, sabemos que

$$E(y) = a_{2l}y^{2l} + a_{2l-2}y^{2l-2} + \dots + a_2y^2 + a_0 = 0 \tag{2.17}$$

con  $l = \max\{j, k\}$  y, como tiene orden  $p$ ,  $E(y) = O(y^{p+1})$ ,  $y \neq 0$ , por lo que tendrá que ser

$$E(y) = a_{2l}y^{2l} + \dots + a_{p+1}y^{p+1} = 0, \quad y \neq 0 \tag{2.18}$$

Como  $R(z)$  es I-estable, entonces si  $y \neq 1$ ,  $E(y) \neq a_{2l}y^{2l} = 0$ , luego  $a_{2l} = 0$ . Por tanto,

$$jQ(iy)^2 - jP(iy)^2 = a_{2l}y^{2l} = 0 \implies jQ(iy)^2 = jP(iy)^2 + a_{2l}y^{2l} = 0$$

$$\implies jR(iy)^2 = \frac{jP(iy)^2}{jQ(iy)^2} = \frac{jP(iy)^2}{jP(iy)^2 + a_{2l}y^{2l}} = 1 \implies jR(1) = 1 \tag{2.19}$$

Tenemos por hipótesis que  $jR(1) = 1$ , entonces necesariamente, el grado del numerador ha de ser menor o igual que el del denominador, esto es,  $k \leq j$ . Por tanto, se obtiene (2.18) con  $l = j$ . Además, cuando  $y \neq 1$ ,  $jR(iy) = 1$  entonces  $E(y) = 0$ ,  $y \neq 1$ . Por otra parte,

$$p = 2j - 2 \implies p + 1 = 2j - 1$$

luego en (2.18) necesariamente tiene que ser

$$E(y) = a_{2j}y^{2j}, \quad \forall y \in \mathbb{R}$$

y  $E(y) = 0$  cuando  $y \neq 1$ . Entonces,  $a_{2j} = 0$  y  $E(y) = 0$ ,  $\forall y \in \mathbb{R}$ . Por lo tanto,  $\mathbb{R}$  es I-estable.

## 2.5. L-estabilidad

**Definición 6** Un metodo se dice L-estable si es A-estable y

$$\lim_{z \rightarrow -\infty} R(z) = 0$$

**Proposición 8** Si un metodo impl cito Runge-Kutta con A no singular satisface una de las siguientes condiciones:

$$\begin{aligned} a_{sj} &= b_j \quad j = 1, \dots, s \\ a_{i1} &= b_1 \quad i = 1, \dots, s \quad (b_1 \neq 0) \end{aligned} \quad (2.20)$$

entonces  $R(-\infty) = 0$ .

*Demostración:* Sabemos por (2.4) que  $R(z)$  es

$$R(z) = 1 + z b^T (I - zA)^{-1} e = 1 + b^T (z^{-1} (I - zA))^{-1} e = 1 + b^T (z^{-1} I - A)^{-1} e$$

Luego, si tomamos  $z \rightarrow -\infty$  nos queda que

$$R(-\infty) = 1 + b^T (-A)^{-1} e = 1 - b^T A^{-1} e$$

Por tanto, si comprobamos cada una de las condiciones (2.20) nos queda:

- Si  $a_{sj} = b_j$ ,  $j = 1, \dots, s$ , y denotamos la última fila de  $A$  como  $A_s^T = e_s^T A$ ,  $e_s^T = (0, \dots, 0, 1)$ , tenemos que  $e_s^T A = b^T$ , y, por lo tanto,

$$R(-\infty) = 1 - b^T A^{-1} e = 1 - e_s^T A A^{-1} e = 1 - e_s^T e = 0$$

- Si  $a_{i1} = b_1$ ,  $i = 1, \dots, s$ ,  $b_1 \neq 0$  vectorialmente podemos denotarlo como que  $Ae_1 = b_1 e$  donde  $e_1 = (1, 0, \dots, 0)^T$ . Por tanto,

$$e = \frac{1}{b_1} Ae_1$$

Por consiguiente,

$$\begin{aligned} R(-\infty) &= 1 - b^T A^{-1} e = 1 - b^T A^{-1} \left( \frac{1}{b_1} Ae_1 \right) = \\ &= 1 - \frac{1}{b_1} b^T A^{-1} Ae_1 = 1 - \frac{1}{b_1} b^T e_1 = 0 \end{aligned}$$

Obsérvese que si  $a_{i1} = b_1 = 0$ ,  $i = 1, \dots, s$ , realmente la primera etapa sería superflua y el método RK sería reducible.

En consecuencia, los métodos A-estables que cumplen alguna de las condiciones (2.20) son L-estables.

### 2.6. Ejemplos

A continuación, aplicaremos lo visto en este capítulo mediante una serie de ejemplos.

(i)  $\theta$ -método [4, p.204]

$$\begin{array}{c|cc}
 0 & 0 & 0 \\
 1 & 1 & \theta \\
 \hline
 1 & \theta & \theta
 \end{array}
 \quad
 R(z) = \frac{1 + z(1 - \theta)}{1 - z\theta}, \quad \theta \in [0, 1] \tag{2.21}$$

Si  $\theta = 0$  tenemos el Euler explícito, si  $\theta = 1$  Euler implícito y si  $\theta = 1/2$  tenemos la regla trapezoidal. En este caso,  $R(z)$  es analítica en  $\mathbb{C}$  dado que para  $\theta > 0$ ,  $z = 1/\theta \notin \mathbb{C}$  y si  $\theta = 0$ ,  $R(z)$  es un polinomio. Así,

$$E(y) = |1 + \theta iy|^2 \quad |1 + (1 - \theta)iy|^2 = (1 + 2\theta)y^2 \quad 0, \quad \theta = \frac{1}{2}$$

Luego, tenemos que el método de Euler implícito y la regla trapezoidal sí son A-estables pero el Euler explícito no lo es. Entre los dos implícitos el único L-estable es el Euler implícito ya que

$$\lim_{z \rightarrow 1} \frac{1 + z(1 - \theta)}{1 - z\theta} = 0, \quad \theta = 1.$$

Obsérvese que  $R(z)$  es la misma para la regla trapezoidal que para la regla del punto medio implícito.

(ii) Método de Hammer-Hollingsworth [4, p.207] de orden 4:

$$\begin{array}{c|ccc}
 \frac{1}{2} & \frac{\rho_3}{6} & & \\
 \frac{1}{2} + \frac{\rho_3}{6} & \frac{1}{4} & \frac{\rho_3}{6} & \\
 \hline
 & & 1/2 & 1/2
 \end{array}
 \quad
 R(z) = \frac{1 + z/2 + z^2/12}{1 - z/2 + z^2/12} \tag{2.22}$$

En este caso,

$$\begin{aligned}
 E(y) &= \left| 1 + \frac{iy}{2} + \frac{(iy)^2}{12} \right|^2 \quad \left| 1 + \frac{iy}{2} + \frac{(iy)^2}{12} \right|^2 = \\
 &= \left| \left( 1 + \frac{y^2}{12} \right) + \frac{iy}{2} \right|^2 \quad \left| \left( 1 + \frac{y^2}{12} \right) + \frac{iy}{2} \right|^2 =
 \end{aligned}$$

$$= \left(1 - \frac{y^2}{12}\right)^2 + \left(\frac{y}{2}\right)^2 - \left(1 - \frac{y^2}{12}\right)^2 - \left(\frac{y}{2}\right)^2 = 0$$

Además, los ceros de  $Q(z)$  son  $z = 3 \pm \sqrt{3}i$ , luego es analítica en  $\mathbb{C}$ , por lo tanto es A-estable. Como  $\lim_{z \rightarrow \infty} R(z) = 1$  no es L-estable.

(iii) SDIRK orden 3 [4, p.207]

$$\begin{array}{c|cc}
 \gamma & \gamma & 0 \\
 1 & \gamma & 1 \\
 \hline
 & 1/2 & 1/2
 \end{array}
 \quad
 R(z) = \frac{1 + z(1 - 2\gamma) + z^2(1/2 - 2\gamma + \gamma^2)}{(1 - \gamma z)^2}
 \quad
 \gamma = \frac{3 \pm \sqrt{3}}{6}$$

(2.23)

En este caso

$$E(y) = (\gamma - 1/2)^2(4\gamma - 1)y^4$$

Por lo tanto, cuando  $\gamma = 1/4$  el método es A-estable ya que además  $R(z)$  es analítica en  $\mathbb{C}$  dado que el cero de  $Q(z)$  es  $z = 1/\gamma > 0$ . Esto implica que será A-estable cuando  $\gamma = (3 + \sqrt{3})/6 \approx 0,7887$  pero no cuando  $\gamma = (3 - \sqrt{3})/6 \approx 0,2113$ . Por otro lado, como  $\lim_{z \rightarrow \infty} R(z) \neq 0$  no es L-estable.

(iv) Método de Kuntzmann y Butcher [4, p.209] de orden 6

$$\begin{array}{c|ccc}
 \frac{1}{2} & \frac{\rho_{15}}{10} & \frac{5}{36} & \frac{2}{9} \\
 \hline
 & \frac{1}{2} & \frac{5}{36} + \frac{\rho_{15}}{24} & \frac{2}{9} \\
 \frac{1}{2} + \frac{\rho_{15}}{10} & \frac{5}{36} + \frac{\rho_{15}}{30} & \frac{2}{9} + \frac{\rho_{15}}{15} & \frac{5}{36} \\
 \hline
 & 5/18 & 4/9 & 5/18
 \end{array}
 \quad
 R(z) = \frac{1 + z/2 + z^2/10 + z^3/120}{1 - z/2 + z^2/10 - z^3/120}$$

(2.24)

En este ejemplo ya vemos que desde que los grados del numerador y denominador aumentan garantizar la A-estabilidad es mucho más complicado. Por eso se desarrolló una teoría más amplia que intenta estudiar cómo deben ser las aproximaciones racionales a  $e^z$  para que sean A-estables, la teoría de las estrellas de orden que veremos en el próximo capítulo.

---

## Estrellas de orden (Order Stars)

Como hemos visto en el capítulo anterior, las funciones de estabilidad  $R(z)$  de los métodos RK son aproximantes racionales a la exponencial, que aproximan su desarrollo de Taylor al menos hasta el orden  $p$  del método. Si queremos obtener aproximaciones a  $e^z$  con orden alto, lo que nos permitirá deducir métodos RK con órdenes altos, los grados  $(k, j)$  de  $R(z)$  tendrán que tener un tamaño relativamente grande, lo que a su vez hace que sea muy difícil saber si es A-estable o no.

Las estrellas de orden van a ayudar a deducir la A-estabilidad de un aproximante racional a  $e^z$  en general y, en particular, determinan la A-estabilidad de la clase de aproximantes con los que se puede alcanzar el mayor orden posible, los aproximantes de Padé.

Hemos seguido principalmente la teoría de las estrellas de orden desarrollada en el [5, Cap.IV.4].

### 3.1. Aproximación de Padé a la función exponencial

Los aproximantes de Padé a la exponencial aportan el orden máximo que se puede obtener con una función racional con grados  $(k, j)$ , del numerador y denominador respectivamente, prefijados.

**Definición 7** [2] *Dada una función  $f$  y dos enteros  $m \geq 0$  y  $n \geq 0$ , el aproximante de Padé de orden  $[m/n]$  es la función racional*

$$R(z) = \frac{a_0 + a_1z + a_2z^2 + \dots + a_mz^m}{1 + b_1z + b_2z^2 + \dots + b_nz^n},$$

que verifica

$$f^{(k)}(0) = R^{(k)}(0), \quad k = 0, 1, \dots, m + n \quad (3.1)$$

Si hacemos el desarrollo en serie de Maclaurin de  $R(z)$ , (3.1) es equivalente a

$$f(z) - R(z) = O(z^{m+n+1}), \quad z \neq 0$$

y es el mayor orden de aproximación posible. Además, el aproximante de Padé es único para  $m$  y  $n$ .

En lo que sigue llamaremos **[k/j]-Padé** al aproximante de Padé de orden  $[k/j]$ .

**Teorema 6** Dado dos enteros  $k \geq 0$  y  $j \geq 0$ , la aproximación  $[k, j]$ -Pade a  $e^z$  esta dada por

$$R_{kj}(z) = \frac{P_k(z)}{Q_j(z)}$$

donde

$$P_k(z) = 1 + \sum_{l=1}^k \frac{(k)_l}{(j+k)_l} \frac{z^l}{l!}, \quad Q_j(z) = 1 + \sum_{l=1}^j (-1)^l \frac{(j)_l}{(k+j)_l} \frac{z^l}{l!} \quad (3.2)$$

denotando  $(k)_l = k(k-1)\dots(k-l+1)$ .

**Ejemplo 1** Veamos algunos ejemplos de aproximantes de Padé.

- $[1/0]$ -Padé:

$$R_{10}(z) = \frac{1+z}{1}$$

Vemos que este aproximante coincide con tiene la función de estabilidad del método de Euler (1.8).

- $[0/1]$ -Padé:

$$R_{01}(z) = \frac{1}{1-z}$$

En este caso, vemos que esta es la función de estabilidad del método de Euler implícito (1.11).

- $[1/1]$ -Padé:

$$R_{11}(z) = \frac{1 + \frac{1}{2}z}{1 - \frac{1}{2}z}$$

Este coincide con la función la función de estabilidad de la regla implícita del punto medio (1.12) y de la regla trapezoidal (1.13).

- [3/3]-Padé

$$R_{33}(z) = \frac{1 + \frac{1}{2}z + \frac{1}{5} \frac{z^2}{2!} + \frac{1}{20} \frac{z^3}{3!}}{1 - \frac{1}{2}z + \frac{1}{5} \frac{z^2}{2!} - \frac{1}{20} \frac{z^3}{3!}}$$

En este caso, coincide con la función de estabilidad del método de Kuntzmann y Butcher de orden 6 (2.24).

### 3.2. Estrellas de orden (Order Stars)

Las estrellas de orden aparecen al comparar un aproximante racional a  $e^z$  con la propia  $e^z$ .

Sea  $R(z) = P(z)/Q(z)$  un aproximante a  $e^z$  de orden  $p \geq 1$  con  $\deg P = k \geq 0$ ,  $\deg Q = j \geq 0$  y coeficientes reales.

**Definición 8** El conjunto

$$A = \{z \in \mathbb{C} : |R(z)| > |e^z|\} = \{z \in \mathbb{C} : |q(z)| > |p(z)|\}, \quad q(z) = R(z)/e^z \quad (3.3)$$

se denomina **estrella de orden (order star)** de  $R$ .

**Proposición 9** La estrella de orden es simétrica respecto al eje real.

*Demostración:* Veamos que  $z \in A$ , entonces  $\bar{z} \in A$ . Al ser  $R(z)$  un cociente de polinomios con coeficientes reales,  $\overline{R(z)} = R(\bar{z})$ ,  $z \in A$ .

Por tanto,  $|R(z)| > |e^z| \iff |\overline{R(z)}| > |\overline{e^z}| \iff |R(\bar{z})| > |e^{\bar{z}}|$ , así

$$|R(\bar{z})| > |e^{\bar{z}}| \iff |e^z| > |R(z)| = |e^{Re z}| = |e^{Re \bar{z}}| = |e^{\bar{z}}|,$$

por lo que  $\bar{z} \in A$ .

**Proposición 10**  $A \setminus i\mathbb{R} = S^c \setminus i\mathbb{R}$ .

*Demostración:* Para todo  $z \in A \setminus i\mathbb{R}$ ,  $z = iy$  para cierto  $y \in \mathbb{R}$  y  $|R(z)| > |e^z| = |e^{iy}| = 1$ ,  $z \notin S$ ,  $z \in S^c \setminus i\mathbb{R}$ .

**Lema 1**  $R(z)$  es I-estable si y solo si

$$A \setminus i\mathbb{R} = \emptyset \quad (3.4)$$

*Demostración:*  $R(z)$  es I-estable si y solo si  $i\mathbb{R} \cap A = \emptyset$ , lo que es equivalente a que  $S^c \setminus i\mathbb{R} = \emptyset = A \setminus i\mathbb{R}$ , por la Proposición 10.

**Lema 2**  $R(z)$  es  $A$ -estable si y solo si  $A \setminus i\mathbb{R} = \emptyset$ ; y todos los polos de  $R(z)$  están en el semiplano positivo  $\mathbb{C}^+$ .

*Demostración:* Sabemos por la Proposición 4 que  $R(z)$  es  $A$ -estable sii es analítica en  $\mathbb{C}^-$  y  $|R(iy)| \leq 1$ ,  $\forall y \in \mathbb{R}$ .

Por tanto, si  $R(z)$  es  $A$ -estable, todos sus polos (o ceros de  $Q(z)$ ) están en  $\mathbb{C}^+$ . Además,  $|R(iy)| \leq 1$ ,  $\forall y \in \mathbb{R}$ ,  $i\mathbb{R} \subset S$  por lo que  $S^c \setminus i\mathbb{R} = \emptyset = A \setminus i\mathbb{R}$  por la proposición 10.

Por otra parte, si todos los polos de  $R(z)$  están en  $\mathbb{C}^+$  entonces  $R(z)$  es analítica en  $\mathbb{C}^-$ . Además,  $A \setminus i\mathbb{R} = S^c \setminus i\mathbb{R} = \emptyset$  implica que  $i\mathbb{R} \subset S$ , es decir,  $|R(iy)| \leq 1$ ,  $\forall y \in \mathbb{R}$ . Por tanto, por la Proposición 4, se tiene la  $A$ -estabilidad.

Conocer cómo son las estrellas de orden asociadas a cada aproximación  $R(z)$  nos dará información para determinar si dicha  $R(z)$  es  $A$ -estable o no. Como estamos tratando con funciones racionales, hay dos casos de  $\mathbb{C}$  en los que es fundamental saber cómo son dichas estrellas: cuando estamos cerca del origen y cuando  $|z| \rightarrow \infty$ .

### 3.3. Order Star cuando $z \rightarrow 0$

Para deducir cómo es  $A$  cerca del origen, es necesario que previamente veamos ciertas propiedades.

**Lema 3** (a) Dada una constante  $0 < k < 1$ ;

$$\cos x > k, \quad x \in D_\gamma^+ := \bigcup_{k \in \mathbb{Z}} (\gamma + 2k\pi, \gamma + 2k\pi), \quad \gamma = \arccos(k) \in \left(0, \frac{\pi}{2}\right].$$

(b) Si  $-1 < k < 0$ ,

$$\cos x < k, \quad x \in D_\gamma := \bigcup_{k \in \mathbb{Z}} (2k\pi + \gamma, (2k+2)\pi - \gamma), \quad \gamma = \arccos(k) \in \left[\frac{\pi}{2}, \pi\right].$$

*Demostración:* Este lema se puede ver fácilmente de forma gráfica. En la figura 3.1 se observa que si elegimos un  $k \in (0, 1)$  la recta  $y = k$  corta a la función  $y = \cos x$  en dos puntos  $\gamma$  y  $-\gamma$  donde  $\gamma = \arccos k$  de tal forma que la función siempre queda por encima de la recta  $y = k$ .

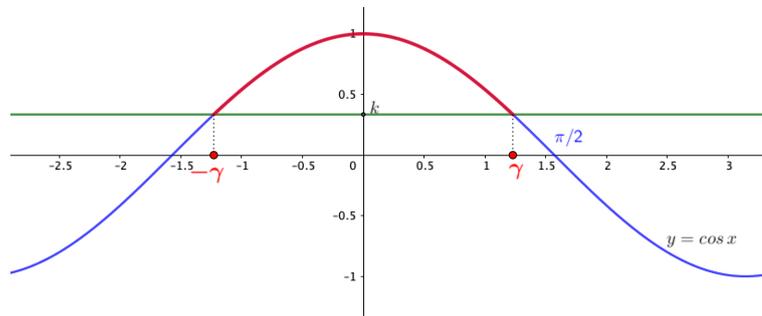


Figura 3.1: Área donde  $\cos x > k$  en rojo

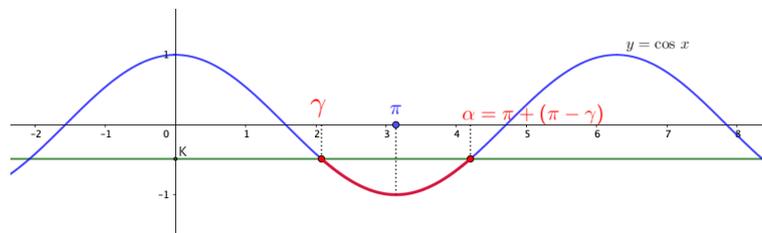


Figura 3.2: Área donde  $\cos x < k$  en rojo

En cambio, si tomamos  $k \in (-1, 0)$  en la Figura 3.2 corta a  $y = \cos x$  en dos puntos  $\gamma = \arccos k \in [\pi/2, \pi]$  y  $\alpha = \pi + (\pi - \gamma)$  de tal forma que el  $\cos x < k$ , para todo  $x \in (\gamma, 2\pi - \gamma)$ . Por la periodicidad del  $\cos x$ , esto que ocurre cuando  $x \in [0, 2\pi]$  se repite en todos los intervalos  $[2k\pi, (2k+2)\pi]$ ,  $k \in \mathbb{Z}$ , de donde se tiene lo que dice el lema.

**Observación 2** De lo anterior se deduce también que:

- Si  $0 < k < 1$  :  $\cos x = k$  ,  $x \in \partial D_\gamma^+$  y  $\cos x < k$  ,  $x \in \mathbb{R} \setminus \overline{(D_\gamma^+)}$
- Si  $-1 < k < 0$  :  $\cos x = k$  ,  $x \in \partial D_\gamma$  ,  $\cos x > k$  ,  $x \in \mathbb{R} \setminus \overline{(D_\gamma)}$

**Lema 4** Dado  $p \in \mathbb{N}$  y una constante  $C \in \mathbb{R}$ ,  $C \neq 0$ , consideramos la función  $f(z) = 1 - Cz^{p+1}$  definida para  $z = re^{i\theta} \in B_M = \{z \in \mathbb{C} / |z| < M\}$  donde se toma  $M > 0$  de forma que  $0 < \frac{|C|}{2} r^{p+1} < 1$  para todo  $0 < r < M$ , en cuyo caso consideramos los siguientes subconjuntos de  $\mathbb{R}$ :

(a) Si  $C > 0$ :  $D_\gamma := \bigcup_{k \in \mathbb{Z}} \left( \frac{\gamma}{p+1} + \frac{2k\pi}{p+1}, \frac{\gamma}{p+1} + \frac{2k\pi}{p+1} \right)$   
 donde  $\gamma = \arccos \left( \frac{C}{2} \right) \in [0, \frac{\pi}{2}]$ .

$$(b) \text{ Si } C < 0: D_\gamma := \bigcup_{k \in \mathbb{Z}} \left( \frac{\gamma}{p+1} + \frac{2k\pi}{p+1}, \frac{2\pi - \gamma}{p+1} + \frac{2k\pi}{p+1} \right)$$

donde  $\gamma = \arccos \left( \frac{Cr^{p+1}}{2} \right) \in \left[ \frac{\pi}{2}, \pi \right]$ .

Entonces, se tiene

- (i) Si  $\theta \in D_\gamma$   $|jf(z)| < 1$
- (ii) Si  $\theta \in \partial D_\gamma$   $|jf(z)| = 1$
- (iii) Si  $\theta \in \mathbb{R} \setminus (D_\gamma \cup \partial D_\gamma)$   $|jf(z)| > 1$

*Demostración:* Primero veamos cómo es el módulo de  $f(z) = 1 - Cr^{p+1}z^{p+1}$ ,  $z = re^{i\theta}$ :

$$\begin{aligned} |f(z)|^2 &= |1 - Cr^{p+1}e^{i(p+1)\theta}|^2 = |1 - Cr^{p+1}(\cos(p+1)\theta + i\sin(p+1)\theta)|^2 = \\ &= (1 - Cr^{p+1}\cos(p+1)\theta)^2 + (Cr^{p+1}\sin(p+1)\theta)^2 = \\ &= 1 - 2Cr^{p+1}\cos(p+1)\theta + C^2r^{2(p+1)}\cos^2(p+1)\theta + C^2r^{2(p+1)}\sin^2(p+1)\theta \end{aligned}$$

Sacando factor común nos queda:

$$|f(z)|^2 = 1 - 2Cr^{p+1}\cos(p+1)\theta + C^2r^{2(p+1)}$$

Por tanto,

$$\begin{aligned} |f(z)| < 1, & \quad 1 - 2Cr^{p+1}\cos(p+1)\theta + C^2r^{2(p+1)} < 1, \\ & \quad , \quad 2Cr^{p+1}\cos(p+1)\theta + C^2r^{2(p+1)} < 0 \end{aligned} \quad (3.5)$$

(a) Si  $C > 0$ , (3.5) es equivalente a

$$Cr^{p+1} - 2r^{p+1}\cos(p+1)\theta < 0 \quad \text{ó} \quad \cos(p+1)\theta > \frac{Cr^{p+1}}{2}$$

Aplicando el apartado (i) del Lema 3 para  $x = (p+1)\theta$  y  $k = Cr^{p+1}/2$ , esto se da si y solo si  $(p+1)\theta \in \left[ \frac{\gamma}{p+1} + \frac{2k\pi}{p+1}, \frac{2\pi - \gamma}{p+1} + \frac{2k\pi}{p+1} \right)$  donde  $\gamma = \arccos(Cr^{p+1}/2)$ , o lo que es lo mismo, si  $\theta \in D_\gamma$ .

De la Observación 3.3 se deducen los puntos (ii) y (iii) del lema.

(b) Si  $C < 0$ , (3.5) es equivalente a que

$$\cos(p+1)\theta < \frac{Cr^{p+1}}{2}$$

y aplicando el apartado (b) del Lema 3 igual que antes se concluye la demostración.

**Lema 5** Suponemos que tenemos las mismas hipótesis del Lema 4 y  $\gamma = \arccos(Cr^{p+1}/2)$ , entonces:

(a) Si  $C > 0$  y  $r \neq 0^+$ , se tiene que  $D_\gamma = D^+$  y

$$D_\gamma \setminus D^+ := \bigcup_{k \in \mathbb{Z}} \left( \frac{\pi/2}{p+1} + \frac{2k\pi}{p+1}, \frac{\pi/2}{p+1} + \frac{2k\pi}{p+1} \right).$$

(b) Si  $C < 0$  y  $r \neq 0^+$ , se tiene que  $D_\gamma = D$  y

$$D_\gamma \setminus D := \bigcup_{k \in \mathbb{Z}} \left( \frac{\pi/2}{p+1} + \frac{2k\pi}{p+1}, \frac{3\pi/2}{p+1} + \frac{2k\pi}{p+1} \right).$$

*Demostración:*

(a) Como  $C > 0$ ,  $r \neq 0^+$ ,  $\gamma = \arccos(Cr^{p+1}/2) \neq \pi/2$  con  $\gamma < \pi/2$  por lo que  $\delta\varepsilon > 0$ ,  $\delta r_0 > 0$  tal que  $\delta r : 0 < r < r_0$ ,  $\gamma = \arccos(Cr^{p+1}/2) < \pi/2$  y  $\pi/2 - \gamma < \varepsilon$ , lo que, aplicando el apartado (a) del Lema 4, demuestra que  $D_\gamma = D^+$  y  $D_\gamma \setminus D^+$ .

(b) Se obtiene de forma similar cuando  $C < 0$  teniendo en cuenta que ahora  $\gamma = \arccos(Cr^{p+1}/2) \neq \pi/2$  con  $\gamma > \pi/2$ .

Obsérvese que, independientemente de que  $C > 0$  o bien  $C < 0$  la recta real se puede dividir como se ve en la Figura 3.3

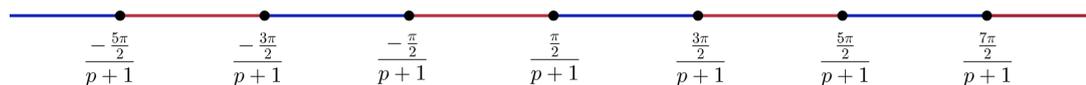


Figura 3.3: Subintervalos de  $D^+$  en rojo y de  $D$  en azul.

**Lema 6** Si  $R(z)$  es una aproximación a  $e^z$  de orden  $p$ , es decir,

$$e^z - R(z) = Cz^{p+1} + O(z^{p+2}), \quad z \neq 0, \quad C \neq 0, \quad (3.6)$$

en un entorno de  $z = 0$  el conjunto  $A$  consiste en los puntos de  $p+1$  sectores circulares de igual amplitud  $\pi/(p+1)$  alternados con  $p+1$  sectores circulares de la misma amplitud de  $A^c$ .

*Demostración:* Dividiendo (3.6) por  $e^z$ , obtenemos

$$\frac{R(z)}{e^z} = 1 - Cz^{p+1} + O(z^{p+2}) = f(z) + O(z^{p+2}), \quad z \neq 0^+ \quad (3.7)$$

donde  $f(z) = 1 - Cz^{p+1}$ . Sea  $z = re^{i\theta}$  con  $r$  suficientemente cerca de 0 para que se verifiquen las hipótesis de los Lemas 4 y 5.

Al igual que en ambos Lemas, tenemos que distinguir si  $C > 0$  o bien  $C < 0$ .

(1) Si  $C > 0$

- Si  $\theta \in D^+$ , por el Lema 5(a) y  $r$  suficientemente pequeño se tiene que  $\theta \in D_\gamma \subset D^+$ . Aplicando el Lema 4(a)-(i), se tiene que  $|f(z)| < 1$ . Por (3.7), se tiene que  $|f(z)| = |R(z)|/e^z < 1$  y, por tanto,  $z \notin A$ .
- Si  $\theta \in D^-$ , entonces  $\theta \in \mathbb{R} \setminus (D_\gamma \cup \partial D_\gamma)$  para todo  $\gamma = \arccos(Cr^{p+1}/2)$ , por lo que, por el Lema 4(a)-(iii) se tiene que  $|f(z)| > 1$  y, por (3.7),  $|f(z)| > 1$  y  $z \in A$ .
- Si  $\theta \in \partial D^+$  no podemos afirmar nada, pero lo anterior nos dice que si consideramos los interiores de los sectores circulares  $fz = re^{i\theta}/\theta \in D^+$   $g$  están contenidos en  $A^c$  y se alternan con los sectores circulares  $fz = re^{i\theta}/\theta \in D^-$   $g$  que están en  $A$ .

(2) Si  $C < 0$ , aplicando los apartados (b) de los Lemas 4 y 5 obtenemos de forma similar los mismo sectores circulares que en el caso  $C > 0$  con la diferencia de que altera el interior de los sectores  $fz = re^{i\theta}/\theta \in D^-$   $g$  están en  $A^c$  y los sectores  $fz = re^{i\theta}/\theta \in D^+$   $g$  están en  $A$ .

Gráficamente



Figura 3.4: Subintervalos de  $D^+$  en rojo y los de  $D^-$  en azul.

**Observación 3** Estudiemos con un poco más de detalle cómo son estos sectores de la estrella de orden cuando  $z \neq 0$ .

Por simetría, el semieje real positivo tiene que ser la bisectriz de un sector y el semieje real negativo la bisectriz de otro.

Por tanto, en el primer cuadrante la primera de las aristas de la “estrella” será la recta

$$\theta = \theta_0 := \frac{\pi}{2(p+1)}$$

Las siguientes aristas serán las rectas

$$\theta = \theta_l := \theta_0 + l \frac{\pi}{p+1}, \quad l = 1, 2, \dots, 2(p+1) - 1.$$

Teniendo en cuenta

$$\theta_0 + l \frac{\pi}{p+1} = \frac{\pi}{2} + \frac{\pi}{2} \left( \frac{1+2l}{p+1} \right) = \frac{\pi}{2} + \frac{1+2l}{p+1} \pi, \quad l = \frac{p}{2}$$

se tiene que

- **Si  $p$  es par:** el semieje imaginario positivo es una arista de la estrella ( $l=p/2$ ) y, por simetría, el semieje imaginario negativo también lo será, por lo que habrá 4 sectores cuyo borde es el eje imaginario. Por tanto, tenemos que hay al menos

$$\frac{p}{2} - 1 \text{ sectores de } A \text{ y } \frac{p}{2} - 1 \text{ sectores de } A^c \text{ completamente contenidos en } C^+. \quad (3.8)$$

y que, al menos,

$$\frac{p}{2} \text{ sectores de } A \text{ y } \frac{p}{2} \text{ sectores de } A^c \text{ arranquen en } C^+. \quad (3.9)$$

- **Si  $p$  es impar:** el semieje imaginario positivo es una bisectriz de un sector de la estrella y el semieje imaginario negativo la bisectriz de otro. Por simetría, ambos pertenecen a  $A$  o ambos pertenecen a  $A^c$ .

Por tanto, tendríamos que al menos hay

$$\frac{p-1}{2} \text{ sectores de } A \text{ y } \frac{p-1}{2} \text{ sectores de } A^c \text{ completamente contenidos en } C^+. \quad (3.10)$$

y, además, hay, como mínimo,

$$\frac{p+1}{2} \text{ sectores de } A \text{ y } \frac{p+1}{2} \text{ sectores de } A^c \text{ que comienzan en } C^+. \quad (3.11)$$

Luego, de forma general tenemos que al menos hay

$$E \left[ \frac{p+1}{2} \right] - 1 \text{ sectores de } A \text{ y } E \left[ \frac{p+1}{2} \right] - 1 \text{ sectores de } A^c \text{ completamente contenidos en } C^+. \quad (3.12)$$

Además hay, al menos,

$$E \left[ \frac{p+1}{2} \right] \text{ sectores de } A \text{ y } E \left[ \frac{p+1}{2} \right] \text{ sectores de } A^c \text{ que comienzan en } C^+. \quad (3.13)$$

Todo lo anterior se da de la misma forma para  $C^-$ .

Por ejemplo, en la Figura 3.5 podemos ver las estrellas de orden del [0/6]-Padé y de la función de estabilidad del método SDIRK de orden 3 (2.6) con  $\gamma = (3 + \sqrt{3})/6$  cerca del origen.

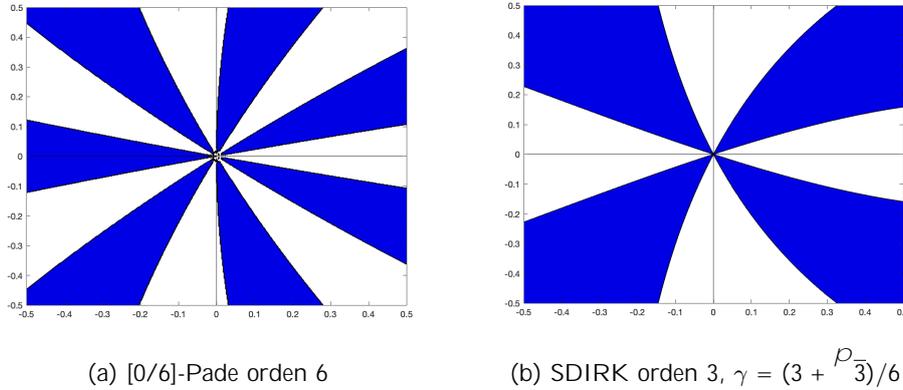


Figura 3.5: Estrellas de orden cerca del origen ( $A$  en azul y  $A^c$  en blanco).

### 3.4. Order Star cuando $z \rightarrow \infty$

Como ya dijimos, otro estudio relevante para entender las estrellas de orden es saber qué ocurre cuando  $z \rightarrow \infty$ .

**Lema 7** Sea  $z = re^{i\theta}$ , entonces

- (i)  $\theta \in \left(\frac{\pi}{2}, \frac{3\pi}{2}\right)$ ,  $\exists K_1 > 0$ ,  $\exists r > K_1$  tal que  $z \in A$ .
- (ii)  $\theta \in \left(-\frac{\pi}{2}, \frac{\pi}{2}\right)$ ,  $\exists K_2 > 0$ ,  $\exists r > K_2$  tal que  $z \notin A$ .

*Demostración:*

- (i)  $\theta \in \left(\frac{\pi}{2}, \frac{3\pi}{2}\right)$ ,  $\cos \theta < 0$ , entonces  $e^{-r \cos \theta} \rightarrow 0, r \rightarrow \infty$ .

Por tanto,

$$\lim_{\substack{r \rightarrow \infty \\ \theta \in (\pi/2, 3\pi/2)}} \frac{|R(z)|}{|e^z|} = \lim_{\substack{r \rightarrow \infty \\ \theta \in (\pi/2, 3\pi/2)}} |R(z)| e^{-r \cos \theta} = +\infty$$

Luego,  $\exists M > 0, \exists K_1 > 0, \exists r > K_1 : \left| \frac{R(z)}{e^z} \right| > M, \theta \in \left(\frac{\pi}{2}, \frac{3\pi}{2}\right)$

Eligiendo  $M > 1, \exists r > K_1, \theta \in \left(\frac{\pi}{2}, \frac{3\pi}{2}\right), z \in A$ .

- (ii)  $\theta \in \left(-\frac{\pi}{2}, \frac{\pi}{2}\right)$ ,  $\cos \theta > 0$ , entonces  $e^{-r \cos \theta} \rightarrow 0, r \rightarrow \infty$

En consecuencia,

$$\lim_{\substack{r \rightarrow \infty \\ \theta \in (-\pi/2, \pi/2)}} \frac{|R(z)|}{|e^z|} = \lim_{\substack{r \rightarrow \infty \\ \theta \in (-\pi/2, \pi/2)}} |R(z)| e^{-r \cos \theta} = 0$$

Luego,  $\exists \varepsilon > 0, \exists K_2 > 0, \exists r > K_2 : \left| \frac{R(z)}{e^z} \right| < \varepsilon, \exists \theta \in \left( \frac{\pi}{2}, \frac{\pi}{2} \right)$   
 Tomando  $\varepsilon < 1, \exists r > K_2, \exists \theta \in \left( \frac{\pi}{2}, \frac{\pi}{2} \right) : z \notin A$ .

**Observación 4** El lema anterior nos garantiza que para todo  $z \in \mathbb{C}$  con  $\operatorname{Im} z > \max\{K_1, K_2\}$  y  $\operatorname{Re} z \neq 0$ ,

$$\operatorname{Re} z < 0, \quad z \in A$$

$$\operatorname{Re} z > 0, \quad z \in A^c.$$

**Lema 8**  $\lim_{z \rightarrow \infty} \left| \frac{R'(z)}{R(z)} \right| = 0$

*Demostración:* Derivando  $R(z)$ :

$$\frac{R'(z)}{R(z)} = \frac{\frac{P'(z)Q(z) - P(z)Q'(z)}{Q(z)^2}}{\frac{P(z)}{Q(z)}} = \frac{P'(z)Q(z) - P(z)Q'(z)}{P(z)Q(z)}.$$

Por tanto, como  $\deg P' = k - 1, \deg Q' = j - 1$ , el numerador tendrá grado  $k + j - 1$  a lo sumo y el denominador grado  $k + j$ . Entonces,  $(P'(z)Q(z) - P(z)Q'(z))/R(z) \rightarrow 0, z \rightarrow \infty$  ya que el grado del numerador es menor que el del denominador.

**Lema 9** Consideramos  $z = re^{i\theta} \in \mathbb{C}$ , con  $r \neq 1$ .  
 Si

$$\varphi_1(\theta) = je^{zj} = e^{2r \cos \theta}, \quad \varphi_2(\theta) = jR(z)j = R(re^{i\theta})R(re^{-i\theta}), \quad (3.14)$$

entonces:

- (i)  $\frac{\varphi_1'(\theta)}{\varphi_1(\theta)} = -2r \sin(\theta)$  y  $\frac{\varphi_2'(\theta)}{\varphi_2(\theta)} = 2r \operatorname{Re}(w)$ , donde  $w = ie^{i\theta} \frac{R'(re^{i\theta})}{R(re^{i\theta})}$ .
- (ii)  $\exists \varepsilon \in (0, \pi/2), \frac{\varphi_1'(\theta)}{\varphi_1(\theta)} < \frac{\varphi_2'(\theta)}{\varphi_2(\theta)}, \exists \theta \in [\varepsilon, \pi - \varepsilon]$ .

*Demostración:*

(i) Derivando cada función, obtenemos:

$$\varphi_1'(\theta) = -2re^{2r \cos(\theta)} \sin(\theta)$$

$$\varphi_2'(\theta) = R'(re^{i\theta}) re^{i\theta} i R(re^{-i\theta}) + R(re^{i\theta}) R'(re^{-i\theta}) re^{-i\theta} (-i)$$

Luego,

$$\frac{\varphi_1^\theta(\theta)}{\varphi_1(\theta)} = \frac{2re^{2r \cos(\theta)} \sin(\theta)}{e^{2r \cos \theta}} = 2r \sin(\theta)$$

$$\begin{aligned} \frac{\varphi_2^\theta(\theta)}{\varphi_2(\theta)} &= \frac{R^\theta(re^{i\theta}) \operatorname{Re}(re^{i\theta}) + R(re^{-i\theta}) R^\theta(re^{-i\theta}) \operatorname{Re}(re^{-i\theta})}{R(re^{i\theta})R(re^{-i\theta})} = \\ &= \frac{R^\theta(re^{i\theta}) \operatorname{Re}(re^{i\theta})}{R(re^{i\theta})} + \frac{R^\theta(re^{-i\theta}) \operatorname{Re}(re^{-i\theta})}{R(re^{-i\theta})} \end{aligned}$$

Dado que los coeficientes de  $R(z)$  son reales, entonces  $R(\bar{z}) = \overline{R(z)}$  y  $R^\theta(\bar{z}) = \overline{R^\theta(z)}$ . Por tanto, si

$$u = \frac{R^\theta(re^{i\theta}) \operatorname{Re}(re^{i\theta})}{R(re^{i\theta})} \quad \bar{u} = \frac{R^\theta(re^{-i\theta}) \operatorname{Re}(re^{-i\theta})}{R(re^{-i\theta})} \quad (3.15)$$

En consecuencia,

$$\frac{\varphi_2^\theta(\theta)}{\varphi_2(\theta)} = u + \bar{u} = 2\operatorname{Re} u = 2r \operatorname{Re}(w), \text{ donde } w = ie^{i\theta} \frac{R^\theta(re^{i\theta})}{R(re^{i\theta})}$$

(ii) Consideremos  $0 < \varepsilon < \pi/2$  y veamos que  $\frac{\varphi_1^\theta(\theta)}{\varphi_1(\theta)} < \frac{\varphi_2^\theta(\theta)}{\varphi_2(\theta)}$ ,  $\forall \theta \in [\varepsilon, \pi - \varepsilon]$ .

Por (i) demostrar (ii) es lo mismo que ver que  $\sin \theta < \operatorname{Re}(w)$ ,  $\forall \theta \in [\varepsilon, \pi - \varepsilon]$ , teniendo en cuenta que  $w = w(r, \theta)$ .

Sabemos por el Lema 8 que si  $z \neq 1$ ,  $\operatorname{Re}(w) \neq 0$ , por tanto,  $\exists \alpha > 0$ ,  $\exists K > 0$ ,  $\forall r > K$ :  $|\operatorname{Re}(w(r, \theta))| < \alpha$ ,  $\forall \theta \in \mathbb{R}$ .

Aplicándolo para  $0 < \alpha = \sin \varepsilon < 1$ , existe  $K > 0$ ,  $\forall r > K$ ,  $|\operatorname{Re}(w(r, \theta))| < \alpha$ ,  $\forall \theta \in [\varepsilon, \pi - \varepsilon]$ .

Además,  $\forall \theta \in [\varepsilon, \pi - \varepsilon]$ :  $\sin \theta > \alpha$ , (ver Figura 3.6). Por tanto,

$$\left. \begin{array}{l} \alpha < \operatorname{Re} w < \alpha \\ \sin \theta < \alpha \end{array} \right\} \Rightarrow \sin \theta < \operatorname{Re} w(r, \theta), \forall \theta \in [\varepsilon, \pi - \varepsilon],$$

con lo que se demuestra (ii).

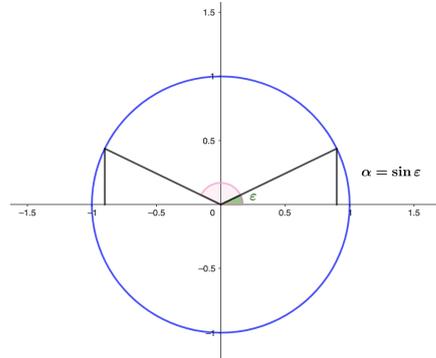


Figura 3.6: Función  $\alpha = \sin \varepsilon$  para  $0 < \varepsilon < \pi/2$ .

**Lema 10** *El borde  $\partial A$  posee solo dos ramas que van al infinito.*

*Demostración:* Para demostrar el lema utilizaremos las funciones definidas en el Lema 9.

Por el Lema 7 sabemos que existen  $K_1, K_2 > 0$  tal que

$$\forall \theta \in (\pi/2, \pi), r > K_1 : z \in A \Rightarrow \varphi_2(\theta) - \varphi_1(\theta) > 0$$

$$\forall \theta \in (0, \pi/2), r > K_2 : z \in A^c \Rightarrow \varphi_2(\theta) - \varphi_1(\theta) < 0$$

Entonces, por continuidad de  $\varphi_1$  y  $\varphi_2$  en  $\mathbb{R}$ ,  $\forall \tilde{\theta} \in (0, \pi) : \varphi_1(\tilde{\theta}) = \varphi_2(\tilde{\theta})$ ,  $\exists r > \max\{K_1, K_2\}$ , lo que implica que  $z = re^{i\theta} \in \partial A$ .

Es decir, hay al menos una rama de  $\partial A$  que se va al infinito en  $Im z > 0$ .

Por el apartado (ii) del Lema 9,  $\forall \varepsilon \in [0, \pi/2], \forall \theta \in [\varepsilon, \pi - \varepsilon]$ :

$$\begin{aligned} \frac{\varphi_1(\theta)}{\varphi_1(\varepsilon)} < \frac{\varphi_2(\theta)}{\varphi_2(\varepsilon)} & \Rightarrow \varphi_1(\theta)\varphi_2(\varepsilon) < \varphi_1(\varepsilon)\varphi_2(\theta) \\ & \Rightarrow \varphi_1(\theta)\varphi_2(\theta) - \varphi_1(\varepsilon)\varphi_2(\theta) < 0 \\ & \Rightarrow \left(\frac{\varphi_1(\theta)}{\varphi_2(\theta)}\right)' = \frac{\varphi_1'(\theta)\varphi_2(\theta) - \varphi_1(\theta)\varphi_2'(\theta)}{\varphi_2^2(\theta)} < 0 \end{aligned}$$

Luego,  $\varphi_1(\theta)/\varphi_2(\theta)$  es estrictamente decreciente en  $[\varepsilon, \pi - \varepsilon]$ ,  $\forall \varepsilon \in (0, \pi/2)$ . Supongamos ahora que hay dos ramas de  $\partial A$  en  $Im z > 0$  que se van a infinito, esto es, existen dos valores  $0 < \theta_1 < \theta_2 < \pi$  tales que  $\varphi_1(\theta_1) = \varphi_2(\theta_1)$ ,  $\varphi_1(\theta_2) = \varphi_2(\theta_2)$ . Entonces, existe  $\varepsilon > 0$  tal que  $0 < \varepsilon < \theta_1 < \theta_2 < \pi - \varepsilon < \pi$  y  $\varphi_1/\varphi_2$  es estrictamente decreciente en  $[\varepsilon, \pi - \varepsilon]$  por lo que

$$\frac{\varphi_1(\theta_1)}{\varphi_2(\theta_1)} > \frac{\varphi_1(\theta_2)}{\varphi_2(\theta_2)}$$

o, lo que es lo mismo, que  $1 > 1$ , lo que es una contradicción que viene de suponer que existían dos valores en los que  $\varphi_1$  y  $\varphi_2$  son iguales. Así que solo existe un valor  $\theta$  en  $[\varepsilon, \pi - \varepsilon]$  tal que  $\varphi_1(\theta) = \varphi_2(\theta)$ .

Consecuentemente, existe una única rama de  $\partial A$  en  $Im z > 0$  que se va al infinito.

Por simetría, habrá otra rama de  $\partial A$  en  $Im z < 0$  que se va al infinito.

**Lema 11** *En las condiciones de los lemas anteriores, si*

$$R(z) = Kz^l + O(z^{l-1}), \quad z \neq 1 \quad (l < 0) \quad (3.16)$$

*estas ramas se acercan asintóticamente a la curva*

$$x = \log |K| + l \log |y| \quad (3.17)$$

*Demostración:* Del Lema 10 tenemos que existe un único  $\tilde{\theta} \in (0, \pi)$  y  $\hat{K} > \max\{K_1, K_2\}g > 0$  tal que  $\delta r > \hat{K} z = re^{i\tilde{\theta}} \in \partial A$ .

Luego, si  $z = x + iy$ , por hipótesis,  $|R(z)| = |e^z| y$

$$|R(z)| = |Kz^l + O(z^{l+1})| = |K|(x^2 + y^2)^{l/2} = |e^z| y$$

$$\Rightarrow \log |K| + \frac{l}{2} \log(x^2 + y^2) = x$$

Como  $z \in \partial A$  podemos acercarnos a  $z$  tanto como queramos con puntos  $A$  y de  $A^c$ , es decir,  $\delta\varepsilon > 0$ ,  $\exists z_1^\varepsilon \in A$ ,  $\exists z_2^\varepsilon \in A^c$  con  $z_1^\varepsilon, z_2^\varepsilon \notin i\mathbb{R}$ , tal que

$$|z_1^\varepsilon| < \varepsilon, \quad \arg z_1^\varepsilon = \tilde{\theta}, \quad |z_1^\varepsilon| > \max\{K_1, K_2\}g$$

$$|z_2^\varepsilon| < \varepsilon, \quad \arg z_2^\varepsilon = \tilde{\theta}, \quad |z_2^\varepsilon| > \max\{K_1, K_2\}g$$

Por la Observación 4,  $Re z_1^\varepsilon < 0$  y  $Re z_2^\varepsilon > 0$ , y se tiene

$$\left. \begin{array}{l} \arg z_1^\varepsilon \in \left(\frac{\pi}{2}, \frac{3\pi}{2}\right) \\ \arg z_2^\varepsilon \in \left(-\frac{\pi}{2}, \frac{\pi}{2}\right) \end{array} \right\} \text{ necesariamente } \tilde{\theta} \neq \frac{\pi}{2},$$

por lo que  $x = r \cos \tilde{\theta} \neq 0$  e  $y = r \sin \tilde{\theta} \neq r$ . En consecuencia,  $x/y \neq 0$ . Por tanto, dado que

$$x^2 + y^2 = y^2(x^2/y^2 + 1) = y^2,$$

las ramas se acercan asintóticamente a

$$x = \log |K| + \frac{l}{2} \log(y^2) = \log |K| + l \log |y|$$

### 3.5. Ceros y polos de $R(z)$

Como  $R(z) = P(z)/Q(z)$  es una función racional y  $e^z \neq 0$ ,  $\forall z \in \mathbb{C}$  se tiene que los ceros de  $q(z)$  son los ceros de  $P(z)$  y los polos de  $q(z)$  son los ceros de  $Q(z)$ . Además,

$$\begin{aligned} \text{Si } w \text{ es cero de } q \text{ entonces, } q(w) = 0 < 1 \text{ por lo tanto, } w \notin A^c. \\ \text{Si } w \text{ es polo de } q \text{ entonces } q(w) = \infty > 1 \text{ por lo que, } w \in A. \end{aligned} \tag{3.18}$$

Combinando las propiedades de las estrellas de orden cerca del origen y cuando  $|z| \rightarrow \infty$  podremos localizar los ceros y polos de la función  $R(z)$ .

**Lema 12** *Cada subconjunto  $F \subset A$  acotado con borde común  $\partial F \subset \partial A$  debe contener al menos un polo de  $q(z)$ .*

*Demostración:* Sea  $F \subset A$  acotado con  $\partial F \subset \partial A$ . Supongamos que no hay polos de  $q(z)$  en  $F$ . Como  $F$  es acotado por el principio del máximo tenemos que

$$\max_{z \in F} |q(z)| = \max_{z \in \partial F \subset \partial A} |q(z)| = 1$$

Pero si  $z \in F \subset A$  entonces  $|q(z)| > 1$ , por lo que, hemos llegado a una contradicción al suponer que no hay polos de  $q(z)$  en  $F$ .

Para estudiar más profundamente el comportamiento de los ceros y polos de  $q(z)$  necesitamos recordar algunos conceptos de análisis complejo, que hemos extraído de [3] y [5].

**Definición 9** *Se dice que  $f$  es **meromorfa** en un abierto  $U$  si es holomorfa en  $U$  salvo en un conjunto de puntos que son polos de  $f$ .*

**Definición 10** *Dadas dos curvas  $\tau_0$  y  $\tau_1 : [0, 1] \rightarrow D$  cerradas y simples en una región  $D$ . Se dice que  $\tau_0$  es **homotópica** a  $\tau_1$  en  $D$  ( $\tau_0 \sim \tau_1$ ) si existe una función continua  $\gamma : [0, 1] \times [0, 1] \rightarrow D$  tal que*

$$\begin{cases} \gamma(s, 0) = \tau_0(s), & \gamma(s, 1) = \tau_1(s), & 0 \leq s \leq 1 \\ \gamma(0, t) = \gamma(1, t), & 0 \leq t \leq 1 \end{cases}$$

Es decir, si  $\tau_1 \sim \tau_2$ , una se puede “deformar” continuamente en la otra.

**Definición 11** *Si  $\tau$  es una curva simple cerrada en  $\mathbb{C}$  y  $a \in \mathbb{C} \setminus \tau$  se define el **índice** de  $\tau$  con respecto al punto  $a$  como*

$$Ind(a, \tau) = \frac{1}{2\pi i} \oint_{\tau} \frac{dw}{w - a}$$

Este valor también se llama **número de giros** de la curva  $\tau$  alrededor del punto  $a$ , ya que si  $\text{Ind}(a, \tau) = n$  entonces  $\tau$  es homotópica a la curva

$$z(t) = a + e^{2\pi i n t}, \quad t \in [0, 1].$$

El Principio del Argumento nos permite contar el número de los ceros y polos de una función meromorfa  $f$  que están rodeados por un contorno simple cerrado  $C$ .

**Teorema 7 (Principio del Argumento).**[3] Sea  $C$  un contorno simple cerrado contenido en un dominio  $D$ . Supongamos que la función  $f$  es meromorfa en  $D$ , y que  $f(z) \neq 0$ ,  $\forall z \in C$ . Entonces,

$$\frac{1}{2\pi i} \oint_C \frac{f'(z)}{f(z)} dz = N_0 - N_p, \quad (3.19)$$

donde  $N_0$  es el número total de ceros de  $f$  y  $N_p$  es el número total de polos de  $f$  en la región de  $D$  rodeada por  $C$ . En la determinación de  $N_0$  y  $N_p$ , tanto los ceros como los polos se cuentan según sus respectivas multiplicidades.

Si hacemos el cambio  $w = f(z)$ ,  $dw = f'(z)dz$  nos queda

$$\frac{1}{2\pi i} \oint_C \frac{f'(z)}{f(z)} dz = \frac{1}{2\pi i} \oint_{f(C)} \frac{dw}{w} = \text{Ind}(0, f(C))$$

Por tanto, el Teorema 9 nos dice que el número de ceros menos el número de polos de  $f$  en la región rodeada por  $C$  es igual al número de giros que la curva  $f(C)$  da alrededor del origen.

**Lema 13** Cada subconjunto acotado  $F \subset A$  con borde común  $\partial F = \partial A$  que reúne  $m$  sectores de  $A$  en el origen debe contener al menos  $m$  polos de  $q(z)$  (contados según su multiplicidad). Análogamente, cada subconjunto acotado  $F^c \subset A^c$  con  $m$  sectores de  $A^c$  en el origen debe contener al menos  $m$  ceros de  $q(z)$ .

*Demostración:* Primero, vamos a demostrarlo en el caso en que  $m = 1$  y  $\partial F$  es una curva cerrada parametrizada orientada positivamente  $c(t)$ ,  $t_0 \leq t \leq t_1$ .

Sea el vector velocidad  $\vec{a} = (c_1'(t), c_2'(t))$ , que es tangente a la curva, y  $\vec{n} = (c_2'(t), -c_1'(t))$  un vector normal exterior a  $\partial F$ .

Escribimos  $q(z) = R(z)/e^z = r(x, y) e^{i\varphi(x, y)}$ , donde  $z = x + iy$ , es decir,  $r(x, y) = |q(z)|$  y  $\varphi(x, y) = \arg(q(z))$ . Luego

$$\log q(z) = \log(r(x, y)) + i\varphi(x, y) \quad (3.20)$$

Por otro lado, para cualquier incremento  $\Delta > 0$ , como el borde  $\partial A$  está rodeado exteriormente por puntos de  $A^c$

$$\begin{cases} (x + \Delta c_2^\ell, y - \Delta c_1^\ell) \in A^c & r(x + \Delta c_2^\ell, y - \Delta c_1^\ell) < 1, \\ (x, y) \in \partial F & \partial A & r(x, y) = 1, \end{cases}$$

entonces,

$$\frac{\partial r}{\partial \vec{n}}(x, y) = \lim_{\Delta \rightarrow 0^+} \frac{r(x + \Delta c_2^\ell, y - \Delta c_1^\ell) - r(x, y)}{\Delta} = 0$$

Por tanto,  $r$  se incrementa hacia dentro de  $F$ , es decir,

$$\frac{\partial r}{\partial \vec{n}}(x, y) \geq 0 \quad \frac{\partial r}{\partial(\vec{n})}(x, y) \leq 0$$

Además,

$$\begin{aligned} \frac{\partial}{\partial \vec{n}}(\log r(x, y)) &= r(\log r) \vec{n} = \left( \frac{r_x}{r}, \frac{r_y}{r} \right) \vec{n} = \\ &= \frac{1}{r} (r_x, r_y) \vec{n} = \frac{1}{r} (r r_x, r r_y) \vec{n} = \frac{1}{r} \frac{\partial r(x, y)}{\partial \vec{n}} = 0 \end{aligned} \quad (3.21)$$

Por otro lado, por las ecuaciones de Cauchy-Riemann para  $\log q(z)$  (3.20),

$$\frac{\partial(\log r)}{\partial x} = \frac{\partial \varphi}{\partial y}, \quad \frac{\partial(\log r)}{\partial y} = -\frac{\partial \varphi}{\partial x} \quad (3.22)$$

de modo que

$$\begin{aligned} \frac{\partial(\log r)}{\partial \vec{n}} &= \left( \frac{\partial(\log r)}{\partial x}, \frac{\partial(\log r)}{\partial y} \right) (c_2^\ell(t), -c_1^\ell(t)) = \frac{\partial(\log r)}{\partial x} c_2^\ell(t) - \frac{\partial(\log r)}{\partial y} c_1^\ell(t) = \\ &= \frac{\partial \varphi}{\partial y} c_2^\ell(t) - \left( -\frac{\partial \varphi}{\partial x} \right) c_1^\ell(t) = \frac{\partial \varphi}{\partial x} c_1^\ell(t) + \frac{\partial \varphi}{\partial y} c_2^\ell(t) = \frac{\partial \varphi}{\partial \vec{a}} \end{aligned}$$

Por tanto, por (3.21),

$$\frac{\partial \varphi}{\partial \vec{a}} = 0. \quad (3.23)$$

Sea  $\tilde{q}(t) = q(c_1(t), c_2(t)) = r(c_1(t), c_2(t)) e^{i\varphi(c_1(t), c_2(t))}$ . Entonces, derivando,

$$\begin{aligned} \tilde{q}'(t) &= \left( \frac{\partial r}{\partial x} c_1'(t) + \frac{\partial r}{\partial y} c_2'(t) \right) e^{i\varphi(c_1(t), c_2(t))} \\ &+ r(c_1(t), c_2(t)) i e^{i\varphi(c_1(t), c_2(t))} \left( \frac{\partial \varphi}{\partial c_1} \frac{\partial c_1}{\partial t} + \frac{\partial \varphi}{\partial c_2} \frac{\partial c_2}{\partial t} \right) = \\ &= \frac{\partial r}{\partial \vec{a}}(c(t)) e^{i\varphi(c(t))} + i q(c(t)) \frac{\partial \varphi}{\partial \vec{a}}(c(t)) \end{aligned}$$

Como  $c(t) \in \partial F$ ,  $\forall t \in [t_0, t_1]$ ,  $r(c(t)) = |q(c(t))| = 1$ ,  $\forall t \in [t_0, t_1]$ , por lo que  $\frac{\partial r}{\partial \vec{a}}(c(t)) = 0$ , se tiene,

$$\tilde{q}^\ell(t) = i q(c(t)) \frac{\partial \varphi}{\partial \bar{a}}(c(t))$$

Por otra parte, por la regla de la cadena,

$$\tilde{q}^\ell(t) = \frac{d}{dt} q(c(t)) = q^\ell(c(t)) c^\ell(t)$$

donde  $c^\ell(t) = c_1^\ell(t) + i c_2^\ell(t)$ . Entonces,

$$i q(c(t)) \frac{\partial \varphi}{\partial \bar{a}}(c(t)) = q^\ell(c(t)) c^\ell(t)$$

Luego,

$$\frac{\partial \varphi}{\partial \bar{a}}(c(t)) = i \frac{q^\ell(c(t))}{q(c(t))} c^\ell(t), \quad \forall t \in [t_0, t_1]$$

Como  $q^\ell$  tiene un número finito de ceros, esto implica que  $\frac{\partial \varphi}{\partial \bar{a}}(c(t)) = 0$  solo en un número finito de puntos y  $\frac{\partial \varphi}{\partial \bar{a}}(c(t)) < 0$  en el resto.

El argumento de  $q(z)$  solo se anula en un número finito de puntos de  $\partial F$  y decrece cuando  $z \in \partial F$  entonces,  $\tau = q(\partial F)$  se recorre en sentido negativo.

Es decir, el argumento de  $q(z)$  decrece (salvo en un número finito de puntos) a medida que se recorre la curva  $\partial F$  en sentido positivo. En otras palabras, la curva  $\tau = q(\partial F)$  se recorre en sentido negativo a medida que  $\partial F$  se recorre en sentido positivo y tiene que dar al menos una vuelta completa.

Además,  $\int_{\tau} \frac{1}{z} dz = 1$ , así que  $\tau = q(\partial F)$  tiene que caer sobre la circunferencia unidad y darle al menos una vuelta en sentido negativo. Por tanto,  $Ind(0, \tau) = -1$ .

Por el Principio del Argumento  $Ind(0, \tau) = N_0 - N_p$  donde  $N_0$  es el número de ceros y  $N_p$  el número de polos de  $q$  en  $F$ . Como  $F \cap A = \emptyset$ , por (3.18), en  $F$  no hay ceros de  $q$ ,  $N_0 = 0$ . Por tanto,  $Ind(0, \tau) = -N_p = -1$ , luego  $N_p = 1$ , así que tiene que haber al menos un polo de  $q$  en  $F$ .

Por otro lado, si  $F \cap A^c$  acotado,  $\partial F \subset \partial A$ , con  $\partial F$  una curva simple cerrada parametrizada  $c(t)$ ,  $t \in [t_0, t_1]$  igual que en el caso anterior, si  $q(z) = r(x, y)e^{i\varphi(x, y)}$  ahora tendremos que

$$\frac{\partial \varphi}{\partial \bar{a}} = \frac{\partial(\log r)}{\partial \bar{a}} = 0$$

y solo se da la igualdad en un número finito de puntos de  $\partial F$ , por lo que  $\varphi$  es creciente, y  $\tau = q(\partial F)$  se recorre en sentido positivo, girando sobre la circunferencia unidad y, por el Principio del Argumento,

$$N_0 - N_p = \text{Ind}(0, \tau)$$

donde  $N_0$  es el número de ceros de  $q$  y  $N_p$  el número de polos en  $F \cap A^c$ . Como  $F \cap A^c$ , no hay polos de  $q$  en  $F$ , entonces  $N_p = 0$  y  $\text{Ind}(0, \tau) = 1$  porque la circunferencia unidad se recorre al menos una vez, luego  $N_0 = 1$ , y  $F$  contiene al menos un cero de  $q$ .

Si  $F \cap A$  con  $\partial F \cap \partial A$  reúne  $m \geq 1$  sectores del origen,

$$F = \bigcup_{i=1}^m F_i \tag{3.24}$$

donde cada uno de los  $F_i$  incluye un solo sector del origen,  $F_i \cap F_j = \emptyset$  y  $\partial F_i \cap \partial A$  es una curva cerrada simple, aplicando lo anterior sobre cada  $F_i$ , se demuestra que hay al menos  $m$  polos de  $q$  en  $F$ .

Si  $F \cap A^c$  es de la forma (3.24) lo que se obtienen son  $m$  ceros de  $q$  en  $F$ .

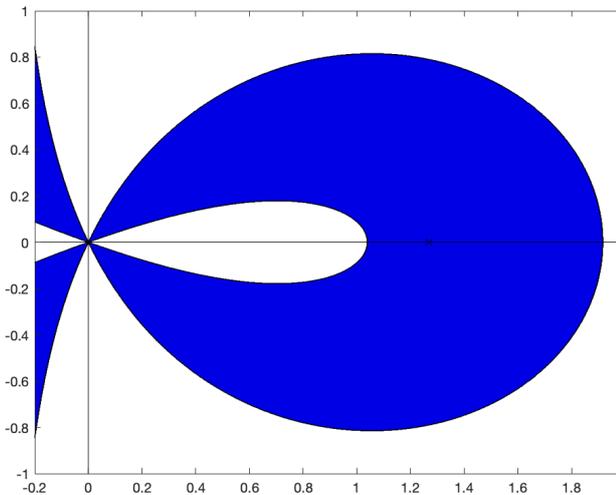


Figura 3.7: Subconjunto  $F \cap A$  del SDIRK  $\gamma = (3 + \sqrt{3})/6$  en color azul.

Si  $F \cap A$  acotado,  $\partial F \cap \partial A$  pero la curva es cerrada no simple, la demostración para un  $m$  general es muy complicada. Veámoslo para el caso  $m = 2$ , es decir, que  $F$  contiene 2 sectores del origen, pues  $\partial F$  no es una curva simple. Este es el caso, por ejemplo, del subconjunto acotado  $F \cap A$  en  $\text{Re } z > 0$  de la estrella de orden de la función de estabilidad del método SDIRK de orden 3 (2.6) con  $\gamma = (3 + \sqrt{3})/6$  que se puede ver en la Figura 3.7.

Así  $\partial F \setminus \partial A$  rodea a un conjunto  $F_2 \setminus A^c$ . Consideremos, además, el subconjunto de  $\mathbb{C}$

$$F_1 = F \setminus F_2$$

Entonces,  $\partial F = \partial F_1 \cup \partial F_2$  y  $\partial F_1$  y  $\partial F_2$  son dos curvas cerradas simples orientadas positivamente.

Sean  $N_0^i$  y  $N_p^i$  el número de ceros en  $F_i$  y  $\tau_i = q(\partial F_i)$  para  $i = 1$  y  $2$ . Sabemos que  $N_0^i - N_p^i = \text{Ind}(0, \tau_i)$ ,  $i = 1, 2$ . Además, al igual que vimos en el caso  $m = 1$ , como  $\partial F_1$  rodea puntos de  $A$ ,  $\tau_1$  recorre al menos una vez la circunferencia unidad en sentido negativo, por lo que  $\text{Ind}(0, \tau_1) = -1$ . En cambio, como  $\partial F_2$  rodea a  $F_2 \setminus A^c$ ,  $\tau_2$  recorre en sentido positivo la circunferencia unidad, al menos una vez al completo, por lo que  $\text{Ind}(0, \tau_2) = 1$ . Como sabemos que en  $F_2 \setminus A^c$  no puede haber polos de  $q$ ,  $N_p^2 = 0$ , por lo que  $N_0^2 = \text{Ind}(0, \tau_2) = 1$ , y si hay polos de  $q$  en  $F$ , han de estar en  $F_1$ .

Además, como en  $F \setminus A$  no puede haber ceros de  $q$ , si  $F_1$  tiene ceros, tienen que estar en  $F_2$ , por tanto,  $N_0^1 = N_0^2$ . Por tanto,

$$N_p^1 = \text{Ind}(0, \tau_2) - \text{Ind}(0, \tau_1) = 1 + 1 = 2$$

con lo que se concluye que debe de haber al menos dos polos de  $F \setminus A$ .

Se procede análogamente si  $F \setminus A^c$  acotado conteniendo 2 sectores del origen.

Para acabar esta sección, hemos creado un programa en Matlab que dibuja las estrellas de orden. En la Figura 3.8 pueden ver algunas estrellas de orden de aproximantes de Padé y de las funciones de estabilidad de los métodos SDIRK de orden 3 que no son Padé.

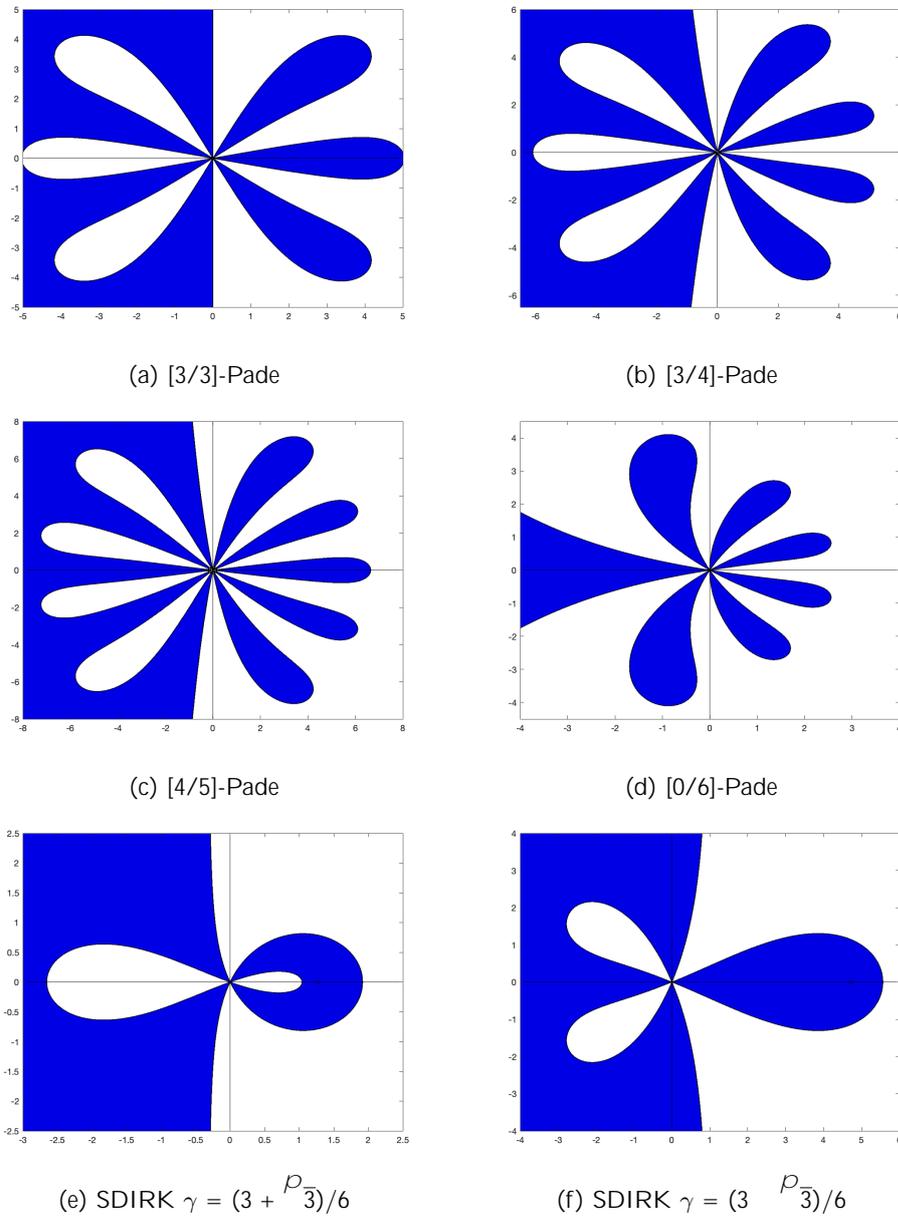


Figura 3.8: Order star

### 3.6. Orden y estabilidad para $R(z)$

Combinando las propiedades vistas anteriormente, podemos deducir algunas condiciones para que una función racional sea A-estable.

**Teorema 8** Si  $R(z)$  es A-estable, entonces  $p \leq 2k_1 + 2$ , donde  $k_1$  es el número diferentes de ceros de  $R(z)$  en  $\mathbb{C}$ .

*Demostración:* Sabemos por la Observación 3 que hay al menos

$$\begin{cases} E \left[ \frac{p+1}{2} \right] & \text{sectores de } A \text{ que arrancan en } \mathbb{C} \\ E \left[ \frac{p+1}{2} \right] - 1 & \text{sectores de } A^c \text{ completamente contenidos en } \mathbb{C} \end{cases}$$

Los sectores de  $A$  que arrancan en  $\mathbb{C}$  no pueden estar acotados, pues si lo estuviera, habría al menos un polo de  $R(z)$  (o de  $q(z)$ ) en  $\mathbb{C}$  por el Lema 13, por lo que no sería  $A$ -estable. Además, los sectores de  $A^c$  completamente contenidos en  $\mathbb{C}$  tienen que estar acotados, pues si fueran infinitos se contradeciría el Lema 7, y cada uno de ellos contendría al menos un cero de  $R(z)$ . Por tanto, habría al menos  $E \left[ \frac{p+1}{2} \right] - 1$  ceros diferentes de  $R(z)$  en  $\mathbb{C}$ , lo que implica que

$$E \left[ \frac{p+1}{2} \right] - 1 \leq k_1$$

En consecuencia,

$$\begin{cases} p - 2k_1 + 1 & \text{si } p \text{ es impar} \\ p - 2k_1 + 2 & \text{si } p \text{ es par} \end{cases}$$

Por tanto, solo podemos garantizar que  $p \geq 2k_1 + 2$ .

**Teorema 9** *Si  $R(z)$  es  $I$ -estable, entonces  $p \geq 2j_1$ , donde  $j_1$  es el número de polos de  $R(z)$  en  $\mathbb{C}^+$ .*

*Demostración:* De igual manera que antes sabemos por la Observación 3 que al menos  $E \left[ \frac{p+1}{2} \right]$  sectores de  $A$  empiezan en  $\mathbb{C}^+$ .

Como  $R(z)$  es  $I$ -estable, los sectores de  $A$  que arrancan en  $\mathbb{C}^+$  no pueden cruzar el eje imaginario, es decir, tienen que estar completamente contenidos en  $\mathbb{C}^+$ , y tienen que estar acotados pues, si no, se contradeciría el Lema 7.

Por el Lema 13, cada uno de esos sectores acotados de  $A$  deben contener al menos un polo de  $R(z)$ , por lo que se garantiza que hay al menos  $E \left[ \frac{p+1}{2} \right]$  polos de  $R(z)$  en  $\mathbb{C}^+$ .

Si  $j_1$  es el número de polos de  $R(z)$  en  $\mathbb{C}^+$ , entonces nos queda

$$E \left[ \frac{p+1}{2} \right] \leq j_1$$

o, lo que es lo mismo,

$$\begin{cases} p \text{ es impar} & p = 2j_1 - 1 \\ p \text{ es par} & p = 2j_1 \end{cases}$$

Entonces, en general solo podemos garantizar que  $p = 2j_1$ .

**Teorema 10** *Supongamos que  $p = 2j - 1$  y  $jR(1)j = 1$ . Entonces,  $R(z)$  es A-estable.*

*Demostración:* Sabemos que  $R(z)$  es A-estable si y solo si, es I-estable y todos los polos están en  $C^+$ .

Como  $p = 2j - 1 > 2j - 2$  entonces es I-estable por la Proposición 7 del Capítulo 2. Por el Teorema 9,  $p = 2j_1 = 2j$  ya que  $j_1 = j$ . Pero además,

$$2j - 1 = p = 2j_1 = 2j - 1 = 2j_1 = j = j_1 + \frac{1}{2}$$

Por tanto,

$$j_1 = j = j_1 + \frac{1}{2} \Rightarrow j = j_1,$$

lo que implica que todos los polos de  $R(z)$  están en  $C^+$ . Aplicando la Proposición 4 del Capítulo 2, se tiene que  $R(z)$  es A-estable.

**Teorema 11** *Suponemos que  $p = 2j - 2$ ,  $jR(1)j = 1$ , y los coeficientes del denominador  $Q(z)$  tienen signos alternos. Entonces,  $R(z)$  es A-estable.*

*Demostración:* Si  $p = 2j - 1$  entonces por el teorema anterior  $R(z)$  es A-estable. Nos falta ver el caso  $p = 2j - 2$ . Por la Proposición 7 del Capítulo 2 sabemos que  $R(z)$  es I-estable. Luego, por el Teorema 9,  $p = 2j_1$  con  $j_1$  el número de polos de  $R(z)$  en  $C^+$ . Por tanto,

$$2j - 2 = p = 2j_1 \Rightarrow j_1 = j - 1$$

entonces nos queda o que  $j_1 = j$  y por tanto A-estable, o bien  $j_1 = j - 1$ . Si  $j_1 = j - 1$  entonces hay un único polo de  $R(z)$  en  $C^-$ , o sea existe  $w \in C^- : Q(w) = 0$ .

Como  $Q(z)$  tiene coeficientes reales,  $w$  tiene que ser real, es decir,  $w = -x \in \mathbb{R}^-, x > 0$ . Sin embargo, por hipótesis

$$Q(z) = b_0 - b_1z + b_2z^2 - b_3z^3 + \dots + (-1)^j b_j z^j$$

con  $b_0 \neq 0$  y  $b_l \neq 0$ . Así que para  $x > 0$ ,

$$\tilde{Q}(x) := Q(x) = b_0 + b_1x + b_2x^2 + \dots + b_jx^j.$$

Podemos ver que no hay cambios de signos en los coeficientes de  $\tilde{Q}(x)$ . Por tanto, no puede tener raíces reales positivas (regla de Descartes) lo que nos lleva a una contradicción, así que no es posible que haya un único polo de  $R(z)$  en  $\mathbb{C}$ , lo que demuestra la A-estabilidad.

Acabamos esta sección con un teorema extraído de [5, p.57], que no demostraremos ya que su demostración requiere más herramientas de análisis complejo de las que no disponemos en un nivel de grado.

**Teorema 12** *Sea  $R(z)$  con  $k_0$  ceros distintos y  $j_0$  polos distintos. Entonces,  $p = k_0 + j_0$ .*

### 3.7. Estabilidad de las aproximaciones de Padé

Aplicando las propiedades de las secciones anteriores podemos deducir la A-estabilidad de los aproximantes de Padé.

**Teorema 13** *La aproximación  $[k/j]$ -Pade, dada en (3.2), es A-estable si y solo si  $k \leq j \leq k + 2$ . Además, todos sus ceros y polos son simples.*

*Demostración:* "  $\Rightarrow$  " Tenemos que ver que  $k \leq j \leq k + 2$ , y tenemos que  $p = k + j$ . Además, por el Teorema 8 sabemos que

$$p = 2k_1 + 2 \leq k + j \leq 2k_1 + 2 = 2k + 2, \quad (k_1 = k).$$

Por tanto,  $j \leq k + 2$  y, como es A-estable, sabemos que  $j \geq k$ .

"  $(\Leftarrow)$  " Tenemos que  $k \leq j \leq k + 2$  y  $p = k + j$ . Por una parte,

$$k \leq j \leq 2 \Rightarrow p = k + j \leq 2j \leq 2k.$$

Por otra parte, por (3.2), si  $j = k$ ,  $jR(1/j) = 1$  y si  $j = k + 1$  ó  $k + 2$ ,  $jR(1/j) = 0$ , así que  $jR(1/j) \neq 1$  para  $k \leq j \leq k + 2$ . Además,  $Q(z)$  tiene signos alternos, así que por el Teorema 11, es A-estable.

Por otro lado, para ver que los polos y ceros son simples tenemos que ver que los ceros y polos tienen multiplicidad 1. Por el Teorema 12 tenemos que si  $k_0$  y  $j_0$  son, respectivamente, el número de ceros y de polos distintos de  $R(z)$ ,

$$k + j = p = k_0 + j_0 \leq k + j \Rightarrow k = k_0 \quad \text{y} \quad j = j_0$$

con lo que se concluye la demostración.

### 3.8. Ejemplos de aproximantes de Padé

En el capítulo 2 vimos que cuando las funciones de estabilidad de los métodos es simple, como en el caso del  $\theta$ -método, detectar si es A-estable o no es sencillo.

Aplicando la teoría de las estrellas de orden desarrollada en este capítulo podemos deducir la A-estabilidad de métodos de alto orden como el caso del método de Kuntzmann y Butcher de orden 6, donde

$$R(z) = \frac{1 + z/2 + z^2/10 + z^3/120}{1 - z/2 + z^2/10 - z^3/120}$$

Como podemos observar la función de estabilidad coincide con el aproximante [3/3]-Padé cuya estrella de orden podemos ver en la Figura 3.8a. Por tanto, podemos aplicar el Teorema 13, y vemos que dicha función es A-estable dado que  $3 \leq 3 \leq 5$ . Por otro lado, como el numerador y el denominador son del mismo grado entonces,  $R(z) \neq 1$  cuando  $z \neq 1$ , por tanto, no es L-estable.

Una aplicación inmediata de esta teoría es que permite construir métodos RK de alto orden A-estable, eligiendo sus coeficientes de forma que su función de estabilidad sea un aproximante de Padé verificando el Teorema 13. Esta es la continuación natural del trabajo realizado en esta memoria.



---

## Conclusiones

Este Trabajo de Fin de Grado se ha centrado en el estudio de los métodos Runge-Kutta, su estabilidad, en particular su A-estabilidad y la teoría de las estrellas de orden.

La estabilidad de los métodos numéricos en la resolución de ecuaciones diferenciales es fundamental para obtener aproximaciones fiables y eficientes. Garantizar que un método es A-estable asegura que propagará concretamente los errores sobre problemas lineales de coeficientes constantes. En la práctica también mejora la estabilidad de los métodos sobre muchos problemas no lineales. Hemos visto que no hay métodos explícitos A-estables por lo que uno de los tópicos de investigación principales en esta área es la de construir métodos RK implícitos A-estables con orden alto.

El principal y último resultado de esta memoria es el Teorema 13, que garantiza que, bajo ciertas condiciones, los aproximantes de Padé producen métodos A-estables con el orden mayor posible. Hemos demostrado las principales propiedades de las estrellas de orden que se usan en la demostración de dicho Teorema.



---

## Bibliografía

- [1] U.M. Ascher and Linda.R. Petzold, *Computer Methods for Ordinary Differential equations and Differential-Algebraic Equations*, SIAM,1997.
- [2] George A. Baker, P. Graves-Morris and P.A. Carruthers *Pade Approximants*, Addison-Wesley Publishing Company, London, 1981.
- [3] John B. Conway, *Functions of One Complex Variable*, Springer-Verlag, New York, 1973.
- [4] E. Hairer, S.P. Nørsett and G. Wanner, *Solving ordinary differential Equations I*, Springer-verlag, 1993.
- [5] E. Hairer and G. Wanner, *Solving ordinary differential Equations II*, Springer-verlag, 1996.
- [6] H. Logemann and E.P. Ryan, *ordinary differential equations*, Springer Verlag, London, 2014.
- [7] G. Söderlind, L. Jay and M. Calvo, *Stiffness 1952-2012: Sixty years in search of a definition*, BIT Number Math, 2015, 55, pp.531-558.



# Order star theory

Yaiza Pérez Tejera

Facultad de Ciencias • Sección de Matemáticas  
Universidad de La Laguna  
alu0101266473@ull.edu.es

## Abstract

In this Bachelor's thesis, the A-stability of Runge-Kutta methods has been studied, with a particular focus on the theory of order stars, which determines necessary and sufficient conditions for the stability function of a method to be A-stable.

To do this, we have divided this work into three chapters. In the first chapter, we introduce the Runge-Kutta methods and some of their properties. In the second chapter, we study their stability and define the stability function and A-stability. Finally, we develop the theory of order stars, which guarantees the A-stability of Padé approximants under certain conditions.

## 1. Introduction to Runge-Kutta Methods

Let us consider the initial value problem

$$y'(t) = f(t, y(t)), \quad y(t_0) = y_0, \quad t \in [t_0, t_f] \quad (1)$$

with  $f : [t_0, t_f] \times \Omega \rightarrow \mathbb{R}^m$ , where  $\Omega \subset \mathbb{R}^m$  is an open set. A solution to IVP (1) is a continuously differentiable function on some interval  $I$  that contains  $t_0$  and it satisfies the IVP.

Given an integer  $s \geq 1$ , the method that, given an approximation  $y_n$  to the solution of the IVP (1) at  $t_n$ , provides an approximation  $y_{n+1}$  to the solution at  $t_{n+1} = t_n + h$  using the formula:

$$k_i = f(t_n + c_i h, y_n + h \sum_{j=1}^s a_{ij} k_j), \quad i = 1, \dots, s \quad (2)$$

$$y_{n+1} = y_n + h \sum_{i=1}^s b_i k_i$$

is called **s-stage Runge-Kutta method**.  $A = (a_{ij})_{i,j=1}^s$  is called the **coefficient matrix**, the vector  $c = (c_1, c_2, \dots, c_s)^T$  is the vector of **nodes**, and  $b = (b_1, b_2, \dots, b_s)^T$  is the vector of **weights**.

**Theorem 1** Let  $f : [t_0, t_f] \times \Omega \rightarrow \mathbb{R}^m$  be a continuous function that satisfies a Lipschitz condition with respect to  $y$ , where  $L$  is the Lipschitz constant. If we take

$$\tilde{h} < \frac{1}{L \max_i \sum_j |a_{ij}|} \quad (3)$$

then for every  $h \in [0, \tilde{h}]$ , there exists a unique solution of the implicit system (2) that can be obtained through functional iteration. If  $f$  is  $p$  times continuously differentiable, the functions  $k_i = k_i(h)$  are also  $C^p([0, \tilde{h}])$ .

**Definition 1** A method  $RK(A, b)$  is of **order  $p \geq 1$**  if for IVPs (1) with  $f \in C^p([t_0, t_f] \times \Omega)$  and Lipschitz, there exists  $K > 0$  independent of  $h$  such that

$$\|y(t_0 + h) - y_1\| \leq K h^{p+1}, \quad h \in [0, \tilde{h}] \quad (4)$$

## 2. Stability of Implicit Runge-Kutta Methods

Applying the RK (2) to the Dahlquist test equation

$$y' = \lambda y, \quad \operatorname{Re} \lambda < 0 \quad (5)$$

it yields

$$y_{n+1} = R(h\lambda)y_n, \quad R(z) = 1 + zb^T(I - zA)^{-1}e \quad (6)$$

where  $e = (1, \dots, 1)^T$  and  $I$  is the identity matrix of size  $s$ . The function  $R(z)$  is called the **stability function** of the method.

## Proposition 1

$$R(z) = \frac{\det(I - zA + zeb^T)}{\det(I - zA)} \quad (7)$$

**Proposition 2** If the method is of order  $p \geq 1$ , then,

$$e^z - R(z) = Cz^{p+1} + O(z^{p+2}) \quad \text{as } z \rightarrow 0, \quad C \neq 0. \quad (8)$$

**Definition 2** A method, whose **stability domain**  $S = \{z \in \mathbb{C} : |R(z)| \leq 1\}$  satisfies  $C^- \cap S \neq \emptyset$  is called **A-stable**.

## 3. Order Stars

Given  $(k, j)$  the  $[k/j]$ -Padé approximant to  $e^z$  is the only rational function that satisfies

$$R_{k,j}(z) = \frac{P_k(z)}{Q_j(z)} = e^z + O(z^{k+j+1}), \quad z \rightarrow 0 \quad (9)$$

with  $\deg P_k = k$ ,  $\deg Q_j = j$ .

**Definition 3** The set

$$A = \{z \in \mathbb{C} : |R(z)| > |e^z| = |z| \in \mathbb{C} : |q(z)| > 1, \quad q(z) = \frac{R(z)}{e^z} \quad (10)$$

is called the **order star** of  $R$ .

**Lemma 1** If  $R(z)$  is an approximation to  $e^z$  of order  $p$ , that is,

$$e^z - R(z) = Cz^{p+1} + O(z^{p+2}), \quad z \rightarrow 0, \quad C \neq 0, \quad (11)$$

In a neighborhood of  $z = 0$ , the set  $A$  consists of the points in  $p+1$  circular sectors of equal width  $\pi/(p+1)$  alternating with  $p+1$  circular sectors of the same width of  $A^c$ .

**Lemma 2** Every bounded subset  $F \subset A$  with a common boundary  $\partial F \subset \partial A$  that includes  $m$  sectors of  $A$  at the origin must contain at least  $m$  poles of  $q(z)$  (counted with their multiplicities). Similarly, every bounded subset  $F \subset A^c$  with  $m$  sectors of  $A^c$  at the origin must contain at least  $m$  zeros of  $q(z)$ .

**Theorem 2** The  $[k/j]$ -Padé approximation, given in (9), is A-stable if and only if  $k \geq j \geq k+2$ . Moreover, all its zeros and poles are simple.

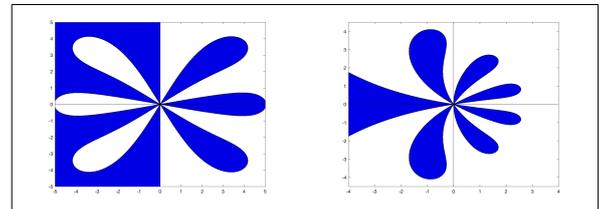


Figure 1: Order star for Padé approximations

## References

- [1] E. Hairer, S.P. Norsett and G. Wanner, "Solving ordinary differential Equations I", Springer-verlag.
- [2] E. Hairer and G. Wanner, "Solving ordinary differential Equations II", Springer-verlag.