

JOILSON B. A. REGO

COMPUTAÇÃO NUMÉRICA

NOTAS DE AULA

 editora
CAULE DE PAPELO®

COMPUTAÇÃO NUMÉRICA

NOTAS DE AULA



COMPUTAÇÃO NUMÉRICA

NOTAS DE AULA

JOILSON B. A. REGO



Natal, 2024





©2024. Joilson B. A. Rego. Reservam-se os direitos e responsabilidades do conteúdo desta edição aos autores. A reprodução de pequenos trechos desta publicação pode ser realizada por qualquer meio, sem a prévia autorização dos autores, desde que citada a fonte. A violação dos direitos do autor (Lei n. 9610/1998) é crime estabelecido pelo artigo 184 do Código Penal. Depósito legal na Biblioteca Nacional conforme Lei N° 10.994, de 14 de dezembro de 2004.

Revisão	<i>O autor</i>
Capa	<i>José Marinho</i>
Projeto Gráfico e Diagramação	<i>Caule de Papiro</i>

Catálogo da Publicação na Fonte.
Bibliotecária/Documentarista:
Rosa Milena dos Santos - CRB 15/847

R343c Rego, Joilson B. A.

Computação numérica: notas de aula [recurso eletrônico] / Joilson B. A. Rego. –
Natal/RN: Caule de Papiro, 2024.

287 p. : il.

ISBN - 978-65-5477-062-0

1. Computação. 2. Cálculo numérico. 3. Matemática aplicada. I. Título.

CDU 681.3.06

Caule de Papiro gráfica e editora
Rua Serra do Mel, 7989, Cidade Satélite
Pitimbu | 59.068-170 | Natal/RN | Brasil
Telefone: 84 3218 4626
www.cauledepapiro.com.br



“O espírito do homem pode descobrir até o infinito, apenas sua preguiça impõe limites à sua sabedoria e suas descobertas.”

(Bossuet)

SUMÁRIO

APRESENTAÇÃO.....	11
-------------------	----

UNIDADE I

CAPÍTULO 1

CONCEITOS BÁSICOS.....	14
Modelagem Matemática e Resolução de Problemas	14
Aproximação	15
Exercícios	18
Propostos	19

CAPÍTULO 2

REPRESENTAÇÃO EM PONTO FLUTUANTE.....	23
Erros	23
Erro absoluto \times erro relativo	25
Sistemas numéricos	25
Conversão de base	27
Parte inteira	27
Conversão de uma base qualquer para uma outra base (qualquer)	30
parte fracionária	33
Representação numérica no computador	34
Representação em ponto fixo	34
Representação em ponto flutuante	35
Operações aritméticas em ponto flutuante	37
Adição e Subtração	37



Multiplicação e divisão em ponto flutuante via método de Horner	38
Erros em soluções numéricas	40
Erro de truncamento	40
Erro de arredondamento	40
Exercícios Propostos	43

CAPÍTULO 3

EXPANSÃO EM SÉRIES DE TAYLOR	47
Polinômio de Taylor	47
Caso Geral	53
Exercícios Resolvidos	62
Exercícios Propostos	69

CAPÍTULO 4

SOLUÇÃO DE EQUAÇÕES NÃO LINEARES	75
Método da Bisseção	78
Regula Falsi ou Método da Falsa Posição	82
Condição de Convergência	84
Ordem de Convergência	84
Método de Newton	85
Método da Secante	88
Iteração funcional e ponto fixo	90
Sistemas de equações não lineares	93
Método de Newton na solução de sistemas não lineares	93
Convergência	97
Exercícios Propostos	98

UNIDADE II

CAPÍTULO 5

SISTEMAS DE EQUAÇÕES LINEARES	105
Métodos diretos	110
Método de Eliminação de Gauss	110
Método de Gauss - Jordan	113



Métodos de Decomposição	114
Decomposição <i>Cholesky</i>	114
Decomposição LU	117
Decomposição LU via método de eliminação de Gauss	118
Matriz inversa com decomposição LU	121
Métodos Iterativos	128
Método iterativo de Jacobi - Richardson	130
Convergência	131
Algoritmo - Método de Jacobi	131
Método de Gauss - Seidel	133
Testes de convergência	133
Algoritmo - Método de Gauss - Seidel	134
Métodos diretos x métodos iterativos	135
Métodos diretos	135
Métodos Iterativos	135
Exercícios Propostos	135

CAPÍTULO 6

AJUSTE DE CURVAS	140
Regressão por Mínimos Quadrados	141
Caso Geral Polinomial	147
Funções de Base	151
Funções de Base Radial - RBF	155
Ajuste não linear	157
Ajuste trigonométrico	160
Exercícios Propostos	167

CAPÍTULO 7

INTERPOLAÇÃO	177
Polinômios interpoladores de Lagrange	178
Polinômios interpoladores por Diferenças Divididas de Newton	187
Exercícios Propostos	



UNIDADE III

CAPÍTULO 8

INTEGRAÇÃO NUMÉRICA	198
Regra do trapézio	200
Regra do trapézio composta	203
Erro de truncamento	204
As Regras de Simpson	205
Limitante superior para o erro	208
Regra 1/3 de Simpson generalizada	208
Limitante superior para o erro	209
Regra 3/8 ou segunda regra de Simpson	210
Quadratura de Gauss	211
Ideia principal	211
Propriedades básicas dos polinômios de Legendre	212
Quadratura de Gauss para $n = 1$	213
Quadratura de Gauss para $n = 2$ e $n = 3$	214
Exercícios Propostos	216

CAPÍTULO 9

RESOLUÇÃO NUMÉRICA DE EDOS	224
Problema de Valor Inicial - PVI	225
Método de Euler	226
Estimativa do erro para o método de Euler	229
Método de Euler modificado	229
Métodos de Runge - Kutta	232
Métodos de Runge-Kutta de Segunda Ordem - RK2	233
Métodos de Runge-Kutta de Terceira Ordem - RK3	237
EDO's de ordem superior	241
Exercícios Propostos	



UNIDADE IV

CAPÍTULO 10

ÁLGEBRA LINEAR	250
Espaços Vetorial	250
Subespaços vetorial	252
Combinação linear	252
Transformações Lineares	254
Matrizes da transformação linear	255
Norma no R^n	257
Produto interno	258
Projeção Ortogonal	261
Matrizes	268
Algumas Matrizes Especiais	269
Matrizes Diagonal e Triangulares	269
Determinante e Matriz Inversa	270
Norma Matricial	276
Autovalores e autovetores	277

CAPÍTULO 11

CÁLCULO DIFERENCIAL	280
Sequências Convergentes	280
Derivadas em Espaços Vetoriais	282
BIBLIOGRAFIA	286



APRESENTAÇÃO

Este trabalho é fruto de uma coletânea das minhas notas de aulas na disciplina de Computação Numérica - CN ministrada nos últimos anos na Escola de Ciência e Tecnologia da Universidade Federal do Rio Grande do Norte (ECT - UFRN). Esta coletânea pretende auxiliar no encaminhamento dos estudos e atividades durante o desenvolvimento da disciplina. Assim, os estudantes podem utilizá-las como material de apoio antes ou após as aulas ministradas.

O principal objetivo ao longo do texto, consiste em apresentar os conceitos essenciais da disciplina e auxiliar na resolução de problemas da Engenharia por meio do entendimento dos métodos estudados para posterior implementação computacional destinada à aferição dos resultados obtidos e apresentar conclusões. Assim, para um bom aproveitamento do curso é recomendável conhecimentos básicos em Álgebra Linear, Cálculo Diferencial e Integral e Linguagem de Programação.

Lembrem-se, como dizia o sábio físico Feynman, o aprendizado e retenção do conhecimento envolve basicamente 04 (quatro) etapas: A primeira etapa consiste na identificação do assunto ou conceito a ser aprendido, em seguida estude-o a fundo, usando os recursos didáticos disponíveis (tais como, estas notas de aula). Após um entendimento básico sobre o conteúdo, procure articular



o que aprendeste (etapa dois). Tente explicar o tópico ou conceitos escolhidos. Nesta etapa, você poderá desmembrar as idéias complexas em partes simples e elementares. A etapa 03 (três) consiste em identificar lacunas (áreas em que travarás ou acabar usando linguagem complexa, porque teu conhecimento não está tão claro e sólido) e voltar ao material de origem para sanar tais lacunas. Finalmente, após entender o assunto bem o suficiente (seja de maneira autônoma ou com ajuda do docente) o conhecimento estará consolidado e permite conexões com o conteúdo em um nível mais aprofundado. Assim, irás perceber a diferença entre “decorar” e compreender verdadeiramente o que foi estudado.

Junho, 2024



UNIDADE I



CONCEITOS BÁSICOS

MODELAGEM MATEMÁTICA E RESOLUÇÃO DE PROBLEMAS

O curso têm como objetivo principal a utilização de ferramentas computacionais aplicadas e eficientes na obtenção de soluções numéricas (aproximadas) onde geralmente não dispomos de uma solução exata ou analítica. A solução numérica pode ser dividida em 05 (cinco) etapas distintas.

- Definição do problema;
- Modelagem Matemática ou formulação do problema;
- Resultados numéricos ou gráficos;
- Implementação computacional;
- Análise e interpretação dos resultados obtidos.

Como exemplo, iremos definir o seguinte problema:

Um paraquedista de massa 78,6 kg pula de um balão de ar quente parado. Deseja-se calcular a velocidade anterior à abertura do paraquedas. Sabendo que o coeficiente de arrasto é igual a 15,5 kg/s. Uma vez definido o problema é necessário (se possível) a modelagem matemática do problema. Ou seja,

$$\sum F = m \cdot a = m \cdot \frac{dv}{dt} \rightarrow F_d - F_v = m \cdot \frac{dv}{dt} \rightarrow g - \frac{c}{m} \cdot v = \frac{dv}{dt}$$



onde, a representa a aceleração, F_d a força devido a aceleração da gravidade, F_v a força devido à resistência do ar, g a aceleração da gravidade, m a massa, c o coeficiente de arrasto e v a velocidade do paraquedista. Portanto, o modelo matemático pode ser definido pela seguinte EDO de primeira ordem:

$$\frac{dv}{dt} = g - \frac{c}{m}v. \quad (1.1)$$

Cuja solução analítica é obtida por meio de um dos métodos de resolução matemática de Equações Diferenciais Ordinárias (EDO's) de primeira ordem. No caso, utilizaremos *fator integrante*.

$$\frac{dv}{dt} + \frac{c}{m}v = g \rightarrow \frac{dv}{dt}e^{(c/m)t} + v\frac{c}{m}e^{(c/m)t} = ge^{(c/m)t} \rightarrow v \cdot e^{(c/m)t} = g \int e^{(c/m)t} dt$$

$$\rightarrow v(t) = g\frac{m}{c} + C \cdot e^{-(c/m)t}$$

considerando $v(0) = 0$, obtemos,

$$v(t) = g\frac{m}{c} \left(1 - e^{-(c/m)t}\right) \quad (1.2)$$

que representa a solução analítica ao problema proposto. No entanto, nem sempre é possível obter uma solução analítica em diversos problemas na Engenharia. Neste caso, utilizaremos o conceito de aproximação.

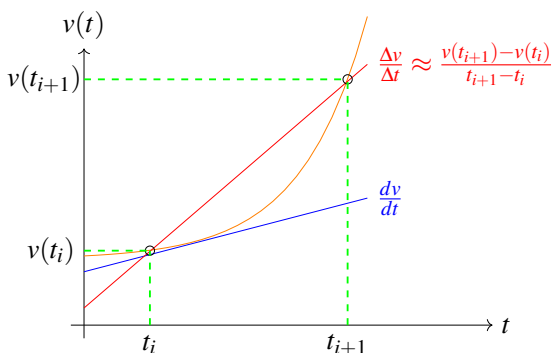
APROXIMAÇÃO

O Principal objetivo é determinar por meio de aproximações, uma função $v(t)$ que satisfaça a EDO, isto é, $v'(t) = f(t, v)$ com $v(t_0) = v_0$ como condição inicial. Assim sendo, iremos aproximar a derivada (reta tangente) pela reta secante num dado intervalo finito fechado. Ou seja,



$$\frac{dv}{dt} \approx \frac{\Delta v}{\Delta t} = \frac{v(t_{i+1}) - v(t_i)}{t_{i+1} - t_i} \quad (1.3)$$

Figura 1.1 - O uso de uma diferença finita para aproximar a derivada da velocidade v com relação ao tempo t



A Equação 1.3 é denominada de *aproximação por diferença dividida finita* da derivada no instante t_i . Substituindo-a na Equação 1.1 obteremos a seguinte relação,

$$v(t_{i+1}) = v(t_i) + \left(g - \frac{c}{m} v(t_i) \right) (t_{i+1} - t_i) \quad (1.4)$$

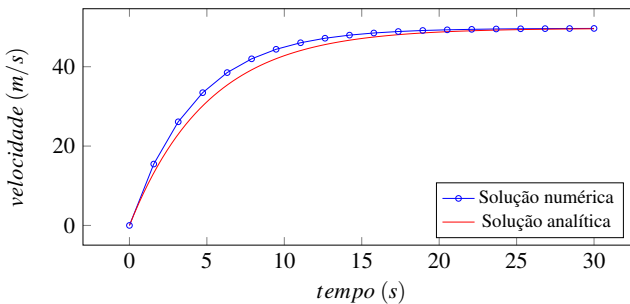
de posse destas informações podemos seguir à próxima etapa que trata da obtenção de resultados numéricos e gráficos para análise dos resultados obtidos através de uma implementação computacional. Considerando inicialmente o paraquedista em repouso, ou seja, $v(t_0) = 0$. Assim, a partir das grandezas definidas no problema inicial, substituindo na equação 1.4, considerando $t_{i+1} = 1s$ para $i = 0$. Obtemos,



$$v(t_1) = v(t_0) + \left(9,8 - \frac{15,5}{78,6}v(t_0)\right)(t_1 - t_0) = 0 + \left(9,8 - \frac{15,5}{78,6} \cdot 0\right)(1 - 0) = 9,8 \text{ m/s}.$$

O resultado está representado graficamente na figura 1.2, onde vemos uma comparação entre a solução analítica e a solução numérica obtida.

Figura 1.2 – Comparação entre as soluções numéricas e analíticas do problema proposto



por último devemos interpretar os resultados tirando algumas conclusões em relação a solução apresentada:

- Claramente surge um erro entre a solução analítica e a numérica, devido ao fato de usarmos aproximações;
- Para minimizar o erro é necessário uma diminuição no passo, ou seja, no Δt ou um aumento no tempo de simulação;
- quando $t \rightarrow \infty$ podemos observar no gráfico da figura 1.2 que a solução numérica tende a solução analítica.

Assim, podemos definir o erro como sendo a diferença entre a solução analítica e a solução numérica. Ou seja, $E = S_a - S_n$. Naturalmente surge um questionamento. Como podemos controlar e identificar o erro



quando não é possível a obtenção da solução analítica? Para isso, iremos abordar em capítulos posteriores o critério de convergência de tais passos.

EXERCÍCIOS

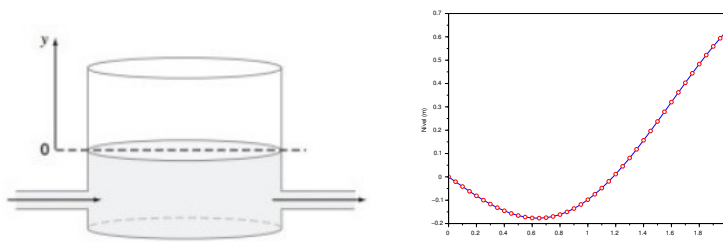
Para fixação do conteúdo apresentado, eis exercícios resolvidos e propostos.

Exercício 1.1

Um tanque de armazenamento contém um líquido à profundidade y , onde $y = 0$ quando o tanque está na metade da capacidade. É tirado líquido a uma vazão constante Q , para atender a demanda, o conteúdo é repostado a uma taxa senoidal de $3Q\sin^2(t)$. A EDO descrita para o sistema é dada por $\frac{d(\Delta y)}{dt} = 3Q\sin^2(t) - Q$. Use o método de aproximação para encontrar a profundidade de y de $t = 0 : 0,5 : 2$. Supondo $A = 1200m^2$, $Q = 500 \frac{m^3}{s}$ e $y(0) = 0$.



Figura 1.3 – Tanque e uma simulação da dinâmica dada pela equação do exercício



$$A \frac{dy}{dt} = 3\text{sen}^2(t) - Q \rightarrow y(t_{i+1}) = \frac{1}{A} (y(t_i) + 3\text{sen}^2(t_i) - Q) (t_{i+1} - t_i)$$

$$= \frac{1}{1200} (y(t_i) + 3\text{sen}^2(t_i) - 500) \cdot (0,5)$$

Tabela 1.1 – Solução do exercício 1.1.

tempo (s)	profundidade (m)
0,0	0,0000
0,5	-0,2083
1,0	-0,2730
1,5	-0,0288
2,0	0,3747

PROPOSTOS

Exercício 1.2

Considere a EDO de primeira ordem: $\frac{dy}{dx} = y \cdot x - x^3$ com $0 \leq x \leq 1,8$ e $y(0) = 1$. Resolva utilizando aproximação com $\Delta x = 0,6$.



Exercício 1.3

A lei do resfriamento de Newton diz que a temperatura de um corpo varia a uma taxa proporcional à diferença entre a sua temperatura e a temperatura do meio que o

cerca (a temperatura ambiente), $\frac{dT}{dt} = -k(T - T_a)$ onde T é a temperatura do corpo (em °C), t é o tempo (min), k é a constante de proporcionalidade (por minuto) e T_a é a temperatura ambiente (°C). Suponha que uma xícara de café originalmente tenha a temperatura de 68°C. Use o método de aproximação da derivada para calcular a temperatura de $t = 0 : 1 : 4$ min se $T_a = 21^\circ\text{C}$ e $k = 0,017$.

Exercício 1.4

Supondo que o corpo de uma vítima foi encontrado em um prédio às 23:00hs, com temperatura corpórea de 30°C. Ao chegar a cena do crime, o investigador solicitou a análise das câmeras de segurança e dados de entrada de pessoas no prédio, obtendo a seguinte tabela,



Tabela 1.2: Exercício 1.4.

Pessoa	horário de trabalho
A	14 - 15:30hs
B	15:45 - 17:30hs
C	17:35 - 18:45hs
D	19 - 19:30hs
E	19:35 - 20:30hs

Sabendo que à meia noite o corpo da vítima estava com temperatura de 29°C e que a temperatura ambiente era de aproximadamente 19°C , sem muita variação. Quem é o principal suspeito?

Exercício 1.5

Algumas células crescem exponencialmente, levando 25 horas para dobrar quando têm um suprimento ilimitado de nutrientes. Entretanto, conforme as células começam a formar uma estrutura esférica e sólida, se houver o corte no suprimento de sangue, o crescimento no centro da estrutura se torna limitado e, eventualmente, as células começam a morrer.

- O crescimento exponencial do número de células N pode ser expresso como mostrado, onde α é a taxa de crescimento das células. Encontre o valor de α .

$$\frac{dN}{dt} = \alpha N$$

- Escreva uma equação que descreva a taxa de variação do volume da estrutura durante o crescimento exponencial, dado que o diâmetro de uma célula individual é de 15 microns.



- Depois que um tipo particular de estrutura passa de 450 microns de diâmetro, as células no centro morrem (mas continuam a tomar espaço). Determine quanto tempo levará para que a estrutura passe deste tamanho crítico.

Exercício 1.6

Mostre que a função $y(x) = x \cdot \text{sen}(x)$ é solução das seguintes EDO's:

- $y + x^2 \cos(x) = xy'$
- $y'' = 2\cos(x) - y$
- $x(y'' + y) = 2y' - 2\text{sen}(x)$

Exercício 1.7

Mostre que a função $y(x) = e^{\text{sen}(x)}$ é solução da EDO:

$$y' = y \cdot \cos(x)$$

considere $y(0) = 1$ e resolva a EDO numericamente por meio de aproximações em seguida, compare a solução numérica com a analítica.



REPRESENTAÇÃO EM PONTO FLUTUANTE

Neste capítulo, abordaremos formas de representar números inteiros e reais em computadores. Iniciaremos com uma discussão sobre a definição e o tipo de erros, tais como erro absoluto e erro relativo, seguido dos sistemas numéricos em diferentes bases e como efetuar a mudança de base entre tais sistemas. Então, enfatizaremos a representação de números com quantidade finita de dígitos, mais especificamente, as representações de números inteiros, ponto fixo e ponto flutuante em computadores.

A representação de números e a aritmética em computadores levam aos chamados erros de arredondamento e de truncamento. Ao final deste capítulo, abordaremos o efeito de tais erros na computação científica.

ERROS

Em geral, problemas na Engenharia são resolvidos de 03 (três) formas distintas: Analiticamente, em que as soluções são baseadas em modelos matemáticos (obtidos das leis físicas dos sistemas), via experimental onde os sistemas são submetidos a ensaios, que representam uma determinada condição de operação e por meio de simulação numérica computacional, em que o modelo é obtido a partir de técnicas em computação numérica. Na maioria dos



casos, inevitavelmente apresentam-se incertezas, sejam aleatórias (inerentes ao problema proposto) ou epistêmicas (que possuem relação com a falta de conhecimento ou experiência). Sendo assim, se possível, por meio de simulações numéricas computacionais, convém verificar se a solução obtida de maneira analítica está convergindo para a solução numérica. Nesta etapa, geralmente se apresentam erros, que pode ser definido por meio de,

$$S_a = S_n + e \rightarrow e = S_a - S_n \quad (2.1)$$

onde S_a representa a solução analítica, S_n a solução numérica e e o erro de aproximação na operação, caracterizado como a diferença entre a solução exata (analítica) e a solução aproximada (numérica). Algumas vezes, não é possível a obtenção de uma solução analítica, neste caso fazemos uso de critérios de convergência para sequências (desde que convergentes) na medição do erro e no critério de parada (abordaremos esta técnica mais adiante).

Exemplo 2.1

$$\text{Ex: } S_a = n! \text{ e } S_n = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n$$

conforme pode ser visto na seguinte tabela.

Tabela 2.1 - Tabela comparativa entre a solução analítica (S_a) × solução numérica (S_n) e o erro apresentado ($S_a - S_n$)

n	S_a	S_n	erro
1	1	0,9221	0,0779
2	2	1,9190	0,0810
3	6	5,8362	0,1638
4	24	23,5062	0,4938
5	120	118,0192	1,9808



ERRO ABSOLUTO × ERRO RELATIVO

Uma vez que definimos o erro como sendo a diferença entre a solução analítica e a numérica, se faz necessário uma melhor acuracidade em relação ao mesmo.

Assim sendo, iremos definir o **erro absoluto** como:

$$E_a = |S_a - S_n| \quad (2.2)$$

caracterizado como a diferença entre a solução analítica e a numérica. E o **erro relativo**,

$$E_r = \frac{E_a}{|S_a|} = \frac{|S_a - S_n|}{|S_a|} (\times 100\%), \quad S_a \neq 0. \quad (2.3)$$

Sendo a relação entre o erro absoluto e a solução analítica. Atualmente, em inúmeros problemas aplicados, existem muitas fontes de erro: tais como, erros simples, de medição, modelagem, etc. No entanto, iremos focar apenas nos erros computacionais, geralmente denominados erros de arredondamento ou truncamento. Para entender tais conceitos, se faz necessário um estudo preliminar de sistemas numéricos e suas representações computacionais.

SISTEMAS NUMÉRICOS

Os sistemas numéricos são sistemas de notação matemática, utilizadas para representar quantidades abstratas denominadas números. Para efeitos práticos, parece bastante razoável adotarmos uma forma de representar tais números de modo a facilitar as operações algébricas envolvidas. Sem, dúvidas qualquer sistema numérico pode ser visto como um polinômio de sua base. De maneira geral,



um número inteiro, no sistema numérico é escrito como sendo uma soma finita, na qual cada parcela é um múltiplo de uma potência da base. Ou seja,

$$p(b) = x_b = \pm (a_n \cdot b^n + a_{n-1} \cdot b^{n-1} + \dots + a_1 \cdot b + a_0) \quad (2.4)$$

onde $x \in \mathbb{Z}$, a_i com $(0 \leq i \leq n)$ são os coeficientes do polinômio, e b representa a base.

Exemplo 2.2

$$x = 354321_{10} = 3 \cdot 10^5 + 5 \cdot 10^4 + 4 \cdot 10^3 + 3 \cdot 10^2 + 2 \cdot 10^1 + 1 \cdot 10^0.$$

No entanto, podemos representar um número em qualquer base, lembrando que os algarismos de representação são sempre de 0 a $b - 1$. Por exemplo, o sistema decimal (base 10), pode ser representado pelos seguintes algarismos (0,1,2,3,4,5,6,7,8,9). O sistema octal (base 8) pode ser representado por (0,1,2,3,4,5,6,7), etc.

Exercício 2.1

Represente os seguintes números inteiros na base correspondente:

- $2743_8 = 2 \cdot 8^3 + 7 \cdot 8^2 + 4 \cdot 8^1 + 3 \cdot 8^0$
- $2743_{10} = 2 \cdot 10^3 + 7 \cdot 10^2 + 4 \cdot 10^1 + 3 \cdot 10^0$
- $21221_3 = 2 \cdot 3^{11} + 1 \cdot 3^{10} + 2 \cdot 3^2 + 2 \cdot 3^1 + 1 \cdot 3^0$, $11_3 = 4_{10}$, $10_3 = 3_{10}$.



CONVERSÃO DE BASE

PARTE INTEIRA

O principal objetivo numa conversão de base, é manipular e entender operações de uma dada representação em uma nova base. Pois, muitos sistemas funcionam com bases distintas, tais como, binária (base 2), octal (base 8) e hexadecimal (base 16). No entanto, para facilitar a plena compreensão, iremos apresentar um algoritmo que nos permite entender e trabalhar com qualquer base. O algoritmo da divisão de Euclides, ou divisão com restos.

Teorema 2.1

O Algoritmo da divisão de Euclides. Dados dois inteiros a e b com $b > 0$, existem únicos inteiros q e r tais que,

$$a = qb + r \tag{2.5}$$

com $0 \leq r < b$, q é denominado de quociente e r o resto da divisão de a por b .

As vezes r também é dito o **resto de a módulo b** . Quando $b > 0$, r é indicado por $r = a \bmod b \equiv \text{mod}(a, b)$.



Exercício 2.2

Converter os seguintes números inteiros na base apresentada para a base solicitada:

$$\left\{ \begin{array}{l} 274 = 91 \cdot 3 + 1 \\ 91 = 30 \cdot 3 + 1 \\ 30 = 10 \cdot 3 + 0 \\ 10 = 3 \cdot 3 + 1 \\ 3 = 1 \cdot 3 + 0 \\ 1 = 0 \cdot 3 + 1 \end{array} \right. \Rightarrow 274_{10} = 101011_3.$$

- $53_{10} = ?_2$

$$\left\{ \begin{array}{l} 53 = 26 \cdot 2 + 1 \\ 26 = 13 \cdot 2 + 0 \\ 13 = 6 \cdot 2 + 1 \\ 6 = 3 \cdot 2 + 0 \\ 3 = 1 \cdot 2 + 1 \\ 1 = 0 \cdot 2 + 1 \end{array} \right. \Rightarrow 53_{10} = 110101_2$$

- $53_{10} = ?_3$

$$\left\{ \begin{array}{l} 53 = 17 \cdot 3 + 2 \\ 17 = 5 \cdot 3 + 2 \\ 5 = 1 \cdot 3 + 2 \\ 1 = 0 \cdot 3 + 1 \end{array} \right. \Rightarrow 53_{10} = 1222_3$$



Uma dúvida natural é, como proceder quando trabalharmos com bases que não é necessariamente a base decimal? (Por exemplo: $23_5 = ?_4$.) Uma maneira prática é utilizarmos a base 10 como ponte, ou seja, proceder da seguinte forma:

$$23_5 = 2 \cdot 5^1 + 3 \cdot 5^0 = 13_{10} \begin{cases} 13 = 3 \cdot 4 + 1 \\ 3 = 0 \cdot 4 + 3 \end{cases} \Rightarrow 23_5 = 31_4$$

Algumas operações elementares entre bases podem ser vistas por meio da seguinte tabela de soma e produto. Por exemplo,

Exemplo 2.3

Construa as tabelas de soma e multiplicação nas bases 3 e 5, respectivamente.

Tabela 2.2 – Tabelas da soma e produto na base 3 e soma na base 5

+	0	1	2		×	0	1	2		+	0	1	2	3	4
0	0	1	2		0	0	0	0		0	0	1	2	3	4
1	1	2	10		1	0	1	2		1	1	2	3	4	10
2	2	10	11		2	0	2	11		2	2	3	4	10	11
										3	3	4	10	11	12
										4	4	10	11	12	13

Tabela 2.3 – Tabela do produto na base 5

+	0	1	2		×	0	1	2		+	0	1	2	3	4
0	0	1	2		0	0	0	0		0	0	1	2	3	4
1	1	2	10		1	0	1	2		1	1	2	3	4	10
2	2	10	11		2	0	2	11		2	2	3	4	10	11
										3	3	4	10	11	12
										4	4	10	11	12	13



CONVERSÃO DE UMA BASE QUALQUER PARA UMA OUTRA BASE (QUALQUER)

Primeiramente iremos efetuar a conversão entre bases utilizando o **Método de Horner** ou **multiplicação alinhada**, que consiste num método eficiente para avaliar polinômios por meio de monômios. O método consiste em reescrever um polinômio de forma a obter uma aproximação para uma certa base.

$$p(b) = \sum_{i=0}^n a_i \cdot b^i = a_0 + a_1 \cdot b + \cdots + a_n \cdot b^n,$$

em que a_0, a_1, \dots, a_n são os coeficientes (reais) do polinômio. No entanto, estamos interessados em calcular seu valor para uma dada base b qualquer. Primeiramente, observe que o polinômio pode ser escrito na forma de parênteses encaixados (ou concatenados):

$$p(b) = a_0 + b(a_1 + b(a_2 + \cdots + b(a_{n-1} + a_n \cdot b) \cdots)).$$

Pelo método podemos definir os números da nova base da seguinte forma:

$$\left\{ \begin{array}{l} x_n = a_n \\ x_{n-1} = a_{n-1} + x_n \cdot b \\ \vdots \\ x_0 = a_0 + x_1 \cdot b \end{array} \right.$$

onde x_0 corresponde o número na nova base.



Exercício 2.3

Converter o número 11201_3 na base 7

$$11201_3 = \underbrace{1}_{a_4} \cdot 3^4 + \underbrace{1}_{a_3} \cdot 3^3 + \underbrace{2}_{a_2} \cdot 3^2 + \underbrace{0}_{a_1} \cdot 3 + \underbrace{1}_{a_0}$$

com

$$\begin{cases} x_4 = a_4 = 1 \\ x_3 = a_3 + x_4 \cdot b = (1 + 1 \cdot 3)_7 = 4_7 \\ x_2 = a_2 + x_3 \cdot b = (2 + \underbrace{4 \cdot 3}_{15_7})_7 = 20_7 \\ x_1 = a_1 + x_2 \cdot b = (0 + 20 \cdot 3)_7 = 60_7 \\ x_0 = a_0 + x_1 \cdot b = (1 + 60 \cdot 3)_7 = 241_7. \end{cases}$$

Portanto, 11201_3 equivale a 241_7 . Ou similarmente, de maneira direta, temos:

$$\begin{aligned} 1201_3 &= 1 \cdot 3^4 + 1 \cdot 3^3 + 2 \cdot 3^2 + 0 \cdot 3 + 1 = 1 + 3(0 + 2 \cdot 3 + 1 \cdot 3^2 + 1 \cdot 3^3) \\ &= 1 + 3(0 + 3(2 + 1 \cdot 3 + 1 \cdot 3^2)) = 1 + 3(0 + 3(2 + 3(1 + 1 \cdot 3))) \\ &= 1 + 3 \left(0 + 3 \left(2 + \underbrace{3(1 + 1 \cdot 3)}_{=4_7} \right) \right) = 1 + 3 \left(0 + 3 \left(\underbrace{2 + 3 \cdot 4}_{=20_7} \right) \right) \\ &= 1 + 3 \left(\underbrace{0 + 3 \cdot 20_7}_{=60_7} \right) = 1 + 3 \cdot 60_7 = 241_7. \end{aligned}$$



Exercício 2.4

Converter o número 372_9 na base 7

$$372_9 = 3 \cdot 9^2 + 7 \cdot 9 + 2 = 3 \cdot (12_7)^2 + 10_7 \cdot (12_7) + 2 = 2 + 12_7 (10_7 + 3 \cdot 12_7)$$

com,

$$\begin{cases} x_2 = 3 \\ x_1 = 10_7 + 3 \cdot 12_7 = 46_7 \\ x_0 = 2 + 46_7 \cdot 12_7 = 620_7. \end{cases}$$

Portanto, 372_9 equivale a 620_7

Um outro método (a partir do método de Horner) foi proposto por Amos O. Olagunju no artigo *A Novel Number conversion Algorithm*. Que pode ser visto através do seguinte exemplo.

Exemplo 2.4

Converter o número 675_8 nas bases 10, 5, 9 e 12.

Tabela 2.4 – Tabela de conversão entre bases utilizando o método proposto por Amos O. Olagunju. Ou seja, $675_8 = 445_{10} = 3240_5 = 544_9 = 311_{12}$

base 8	base 10	base 5	base 9
6	$0 \times 8 + 6 = 6$	$0 \times 13 + 11 = 11$	$0 \times 8 + 6 = 6$
7	$6 \times 8 + 7 = 55$	$11 \times 13 + 12 = 210$	$((6 \times 8)_9 + 7)_9 = 61$
5	$55 \times 8 + 5 = 445$	$210 \times 13 + 10 = 3240$	$((61 \times 8)_9 + 5)_9 = 544$

base 8	base 12
6	$0 \times 12 + 6 = 6$
7	$((6 \times 8)_{12} + 7)_{12} = 47$
5	$((47 \times 8)_{12} + 5)_{12} = 311$



PARTE FRACIONÁRIA

Na parte fracionária utilizaremos um método de multiplicação sucessiva pela base.

Exemplo 2.5

Converter os seguintes números fracionários na base apresentada para a base solicitada:

- $0,7_{10} = ?_3$

$$\left\{ \begin{array}{l} 0,7 \cdot 3 = 0,1 + 2 \\ 0,1 \cdot 3 = 0,3 + 0 \\ 0,3 \cdot 3 = 0,9 + 0 \\ 0,9 \cdot 3 = 0,7 + 2 \Rightarrow 0,7_{10} = 0,2002200220 \dots_3 \\ 0,7 \cdot 3 = 0,1 + 2 \\ 0,1 \cdot 3 = 0,3 + 0 \\ \vdots \end{array} \right.$$

- $0,7_{10} = ?_2$

$$\left\{ \begin{array}{l} 0,7 \cdot 2 = 0,4 + 1 \\ 0,4 \cdot 2 = 0,8 + 0 \\ 0,8 \cdot 2 = 0,6 + 1 \\ 0,6 \cdot 2 = 0,2 + 1 \\ 0,2 \cdot 2 = 0,4 + 0 \Rightarrow 0,7_{10} = 0,101100110 \dots_2 \\ 0,4 \cdot 2 = 0,8 + 0 \\ 0,8 \cdot 2 = 0,6 + 1 \\ \vdots \end{array} \right.$$



Exemplo 2.6

Resolva a seguinte equação e apresente o resultado em binário.

$$(2^5 - 2^4)x = 1011_2$$

$$(2^5 - 2^4)x = 1011_2 \rightarrow 2^4(2-1)x = 1.2^3 + 0.2^2 + 1.2^1 + 1.2^0$$

$$2^4 \cdot 2^{-4}(2-1)x = (1.2^3 + 0.2^2 + 1.2^1 + 1.2^0) \cdot 2^{-4} \rightarrow x = 2^{-1} + 2^{-3} + 2^{-4} \rightarrow x = 0,6875_{10}$$

$$\begin{cases} 0,6875 \cdot 2 = 0,375 + 1 \\ 0,3750 \cdot 2 = 0,750 + 0 \\ 0,7500 \cdot 2 = 0,500 + 1 \\ 0,5000 \cdot 2 = 0,000 + 1 \end{cases} \rightarrow x = 0,6875_{10} = 0,1011_2$$

REPRESENTAÇÃO NUMÉRICA NO COMPUTADOR

Nesta seção, apresentamos um modelo para aritmética computacional de números em ponto flutuante. Existem vários modelos, mas para simplificar, escolheremos um modelo específico e o descreveremos. O formato aritmético em ponto flutuante tornou-se o padrão comum para aritmética de precisão simples e dupla em todo a indústria de computadores.

REPRESENTAÇÃO EM PONTO FIXO

O sistema em ponto fixo representa as partes inteiras e fracionárias de um dado número com uma quantidade fixa de dígitos. Portanto, para uma determinada notação em ponto fixo, indica-se apenas a quantidade de dígitos alocados para a parte inteira e a fracionária. Sendo assim, todos os números a serem manipulados seguem a



mesma notação (fixa). Um número representado em ponto fixo, segue a seguinte notação:

$$x = \pm \sum_{k=m}^n x_k b^{-k}$$

onde:

- $m, n \in \mathbb{Z}$
- $0 \leq x_i \leq b-1$
- $m \leq 0$ e $n > 0$

REPRESENTAÇÃO EM PONTO FLUTUANTE

Basicamente a representação em ponto flutuante consiste de três partes distintas: o sinal (+ ou -), uma mantissa, que contém uma sequência de bits significativos e um expoente.

As três partes são armazenadas juntas em um única palavra, geralmente num computador. Basicamente um sistema de numeração em ponto flutuante pode ser visto como:

$$F(b, p, m, M) = \pm 0, d_1 d_2 \dots d_p \cdot b^e$$

em que,

- b representa a base;
- p a precisão da máquina, ou seja o número de bits da mantissa;
- $m \leq e \leq M$ é o expoente, onde m representa o menor e M o maior expoente;
- obrigatoriamente $d_1 \neq 0$.



A precisão p de uma máquina representa a quantidade de dígitos significativos utilizados na representação numérica. Num sistema em ponto flutuante o menor número real positivo que pode ser representado no sistema $F(b, p, m, M)$ é,

$$x_m = 0, \underbrace{100 \cdots 0}_{p \text{ dígitos}} \cdot b^m = b^{-1} \cdot b^m = b^{m-1}$$

onde m representa o menor expoente e, o maior real positivo representado no mesmo sistema é dado por,

$$x_M = 0, \underbrace{(b-1)(b-1) \cdots (b-1)}_{p \text{ dígitos}} \cdot b^M = (1 - b^{-p}) \cdot b^M$$

Exemplo 2.7

Considere o sistema $F(2, 3, -1, 2)$. Quantos números podem ser representados neste sistema?

$$\begin{pmatrix} 0, 100 \cdot 2^{-1} \\ 0, 101 \cdot 2^{-1} \\ 0, 110 \cdot 2^{-1} \\ 0, 111 \cdot 2^{-1} \end{pmatrix} \begin{pmatrix} 0, 100 \cdot 2^0 \\ 0, 101 \cdot 2^0 \\ 0, 110 \cdot 2^0 \\ 0, 111 \cdot 2^0 \end{pmatrix} \begin{pmatrix} 0, 100 \cdot 2^1 \\ 0, 101 \cdot 2^1 \\ 0, 110 \cdot 2^1 \\ 0, 111 \cdot 2^1 \end{pmatrix} \begin{pmatrix} 0, 100 \cdot 2^2 \\ 0, 101 \cdot 2^2 \\ 0, 110 \cdot 2^2 \\ 0, 111 \cdot 2^2 \end{pmatrix}$$

Exemplo 2.8

Qual o menor e o maior número real positivo representado no sistema $F(2, 3, -3, 3)$?

$$x_m = 2^{(-3-1)} = 2^{-4} = \left(\frac{1}{16}\right)_{10}$$

$$x_M = (1 - 2^{-3}) \cdot 2^3 = 2^3 - 2^0 = 7_{10}$$



OPERAÇÕES ARITMÉTICAS EM PONTO FLUTUANTE

As operações em pontos flutuantes são encontradas normalmente em sistemas que operam em uma grande faixa, de número muito grandes ou muito pequenos, que exige tempo de processamento rápido. Um número, em geral, é representado para um número fixo de dígitos significativos, e escalado usando um expoente em alguma base fixa. Nos sistemas digitais, a base utilizada é a base binária. Desde a década de 1990, a representação mais usada é a definida pelo padrão IEE 754.

ADIÇÃO E SUBTRAÇÃO

As operações de adição e subtração seguem algumas regras pré-definidas. Assim, sejam dois números x e y . Representados no formato $F(b, p, m, M)$. Temos os seguintes passos,

1. Escolher o número com menor expoente, entre os dois, e deslocar sua mantissa à direita um total de dígitos igual a diferença absoluta entre os respectivos expoentes;
2. Colocar o expoente do resultado igual ao maior expoente entre os números;
3. Executar a adição ou subtração das mantissas e determinar o sinal do resultado;
4. Se necessário, normalizar o resultado;
5. Se necessário, fazer arredondamento;
6. Verificar se houver under ou overflow.



Exemplo 2.9

Sejam $x = 0,4546 \times 10^5$ e $y = 0,5433 \times 10^7$, no sistema decimal. Calcule $z = x + y$ utilizando aritmética em ponto flutuante.

Como os expoentes são diferentes, no número de menor expoente deve ocorrer um deslocamento (à direita) na mantissa da diferença entre os expoentes $7 - 5 = 2$. Ou seja, $x = 0,0045 \times 10^7$.

Desde modo,

$$z = 0,0045 \times 10^7 + 0,5433 \times 10^7 = (0,0045 + 0,5433) \times 10^7 = 0,5478 \times 10^7 \text{ (normalizado)}.$$

Exemplo 2.10

Sejam $x = 2,25$ e $y = 134,0625$, no sistema decimal. Calcule $z = x + y$ utilizando aritmética em ponto flutuante.

Sejam,

$$x = 2,25 \times 10^0 = 0,0225 \times 10^2 \text{ e } y = 1,340625 \times 10^2$$

que correspondem as operações dos passos 1 e 2. Com isso,

$$z = 0,0225 \times 10^2 + 1,340625 \times 10^2 = (0,0225 + 1,340625) \times 10^2 = 1,363125 \times 10^2 \text{ (normalizado)}.$$

MULTIPLICAÇÃO E DIVISÃO EM PONTO FLUTUANTE VIA MÉTODO DE HORNER

O método de Horner requer que o multiplicador ou o divisor sejam conhecidos antecipadamente. Isso não constitui uma limitação, tendo em vista que poucos aplicativos realizam multiplicação ou divisão de números que mudam em tempo de execução. Depois que esse valor é estabelecido, a multiplicação ou divisão pode ser realizada de forma eficiente com apenas deslocamentos e operações



simples. O operando é denotado por X , o multiplicador por M e o divisor por D .

- Passo 01 : $X_1 = X \times 2^{-2} + X$.
- Passo 02 : $X_2 = X_1 \times 2^{-3} + X$.
- Passo 03 : $X_3 = X_2 \times 2^{-3} + X$.
- Passo 04 : *resultado final* = $x_3 \times 2^{-3}$

Exemplo 2.11

Realize a seguinte operação na base 2 $(0, 12345_{10}) \cdot (0, 14325_{10})$. Primeiramente precisamos converter os números para a base binária. Assim, temos: $X = 0, 12345_{10} = 0, 000111111001_2$ e $M = 0, 14325_{10} = 0, 001001001010_2$ na representação binária com 12 bits. Em seguida, obtemos

$$\text{Passo 01 : } X_1 = X \times 2^{-2} + X$$

$$\begin{aligned} X_1 &= (0.2^{-1} + 0.2^{-2} + 0.2^{-3} + 1.2^{-4} \dots + 1.2^{-12})2^{-2} + 0,000111111001 \\ &= 0,000001111110 + 0,000111111001 = 0,001001110111. \end{aligned}$$

$$\text{Passo 02 : } X_2 = X_1 \times 2^{-3} + X$$

$$\begin{aligned} X_2 &= (0.2^{-1} + 0.2^{-2} + 1.2^{-3} + 0.2^{-4} \dots + 1.2^{-12})2^{-3} + 0,000111111001 \\ &= 0,0000001001110 + 0,000111111001 = 0,001001000111. \end{aligned}$$

$$\text{Passo 03 : } X_3 = X_2 \times 2^{-3} + X$$

$$\begin{aligned} X_3 &= (0.2^{-1} + 0.2^{-2} + 1.2^{-3} + 0.2^{-4} \dots + 1.2^{-12})2^{-3} + 0,000111111001 \\ &= 0,000001001000 + 0,000111111001 = 0,001001000001. \end{aligned}$$

$$\text{Passo 04 : } X_f = X_3 \times 2^{-3}$$

$$X_f = (0,001001000001)2^{-3} = 0.000001001000_2 = 0,017578125_{10}.$$

com um erro absoluto de aproximadamente 10^{-4} .



ERROS EM SOLUÇÕES NUMÉRICAS

O número excessivo de operações matemáticas executadas na resolução de um dado problema, pode introduzir alguns erros de arredondamento ou dependendo da precisão da máquina, erros de truncamento.

A eficácia do método computacional utilizado na resolução é diretamente afetada pelos efeitos dos erros de truncamento e arredondamento. Em alguns casos, embora os valores exatos sejam supostamente conhecidos, não podem ser exatamente representados na máquina, devido a certas limitações de hardware. Sendo assim, podemos ver algumas categorias de erros que poderemos trabalhar, tais como:

ERRO DE TRUNCAMENTO

Truncar significa descartar todos os dígitos menos significativos a partir de uma determinada posição.

Exemplo 2.12

Truncar o número $1, 01110_2$ para 3 dígitos significativos.

$$1,01_2 = 0,101_2^1$$

ERRO DE ARREDONDAMENTO

Arredondamento é o processo de escolha da representação de um número real em um sistema numérico de ponto flutuante. Para um determinado sistema numérico e um procedimento de arredondamento, o épsilon de máquina é o máximo erro relativo do procedimento escolhido. De maneira prática, arredondar



significa truncar todos os dígitos menos significativos a partir de uma determinada posição, de forma a torná-lo mais próximo do número original.

Regra para arredondamento em representação binária:

$$x = b_1 b_2 \cdots b_d b_{d+1} \cdots b_i \quad b_i \in \{0, 1\}.$$

Trunca-se x na posição d : $x' = b_1 b_2 \cdots b_d$

- se $b_{d+1} = 0 \rightarrow x'$
- se $b_{d+1} = 1 \rightarrow x' = b_1 b_2 \cdots (b_d + 1)$

Exemplo 2.13

Represente o número 1, 45_{10} no sistema $F(2, 3, -3, 3)$ utilizando truncamento, arredondamento e calcule o erro relativo associado.

$$\left\{ \begin{array}{l} 0,45 \cdot 2 = 0,9 + 0 \\ 0,90 \cdot 2 = 0,8 + 1 \\ 0,80 \cdot 2 = 0,6 + 1 \\ 0,60 \cdot 2 = 0,2 + 1 \rightarrow 0,45_{10} = 0,011100110 \cdots_2 \\ 0,20 \cdot 2 = 0,4 + 0 \\ 0,40 \cdot 2 = 0,8 + 0 \\ \vdots \end{array} \right.$$

$$1,45 = 1,011100110 \cdots_2 \rightarrow 0,1011100110 \cdots \cdot 2^1$$

com truncamento temos,

$$1,45 = 1,011100110 \cdots_2 \rightarrow 0,101 \cdot 2^1 = (1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3}) 2^1 = 2^0 + 2^{-2} = 1,25_{10}$$

e erro relativo $\epsilon_r = \frac{|1,45 - 1,25|}{1,45} = 0,1379 = 13,79\%$

com arredondamento temos,



$$\left\{ \begin{array}{l} 0,101 \\ + 1 \end{array} \right. = 0,110 = (1.2^{-1} + 1.2^{-2})2^1 = 2^0 + 2^{-1} = 1,5_{10}$$

erro relativo,

$$E_r = \frac{|1,45 - 1,5|}{1,45} = 0,0344 = 3,44\%.$$

Cabe aqui apresentar um conceito extremamente importante em aritmética de ponto flutuante.

Definição 2.1

Épsilon de máquina.

Denomina-se épsilon da máquina (ϵ) o menor número que somado a 1 produza resultado diferente de 1, ou seja, que não é arredondado. O épsilon de máquina representa a exatidão relativa da aritmética do computador, e a sua existência é uma consequência direta da precisão finita da aritmética de ponto flutuante. Geralmente calculado por meio das seguintes expressões,

$$\epsilon = b^{-(p-1)}$$

Algumas explicações são necessárias para se determinar o valor dessa definição. Um sistema numérico de ponto flutuante é caracterizado por uma base b , e por uma precisão p , por exemplo, o número de dígitos na base b da mantissa (incluindo qualquer bit implícito). Todos os números com o mesmo expoente e possuem espaçamento $b^{e-(p-1)}$. O espaçamento muda nos números que são potências perfeitas de b ; o espaçamento no lado de maior magnitude é b vezes maior que o espaçamento no lado de menor magnitude.



Uma vez que o ϵ de máquina é o limite do erro relativo, é suficiente considerar números com expoente $e = 0$.

A computação numérica usa o ϵ de máquina para estudar os efeitos dos erros de arredondamento. O padrão aritmético da IEEE define que todas as operações de ponto flutuante são feitas como se fosse possível executá-las com precisão infinita, e então o resultado é arredondado para um número de ponto flutuante.

EXERCÍCIOS PROPOSTOS

Exercício 2.5

Escreva os seguintes números binários na notação de ponto flutuante e em seguida converta-os para decimal:

- $11000101, 101_2$
- $110100, 001111_2$

Exercício 2.6

Liste todos os números em ponto flutuante que são expressos na forma:

$$x = \pm(0, d_1 d_2 d_3)_2 \times 2^{\pm k}, k \in \{0, 1\}$$



Exercício 2.7

Considere o sistema em ponto flutuante $F(2, 3, -2, 2)$.
Para esse sistema,

- qual o menor número positivo representável na base 10?
- qual o maior número positivo representável na base 10?

Exercício 2.8

Construa as tabelas de adição e multiplicação para as seguintes bases:

- base 7
- base 8

Exercício 2.9

Calcule os valores obtidos nas seguintes operações:

- $11112+10112 =$
- $89A12+5A612 =$
- $10768+20768 =$
- $307616+577616 =$
- $89A12+89A16 =$



Exercício 2.10

Represente o valor de $f'(2)$ para $f(x) = x^2$ utilizando a aproximação

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}$$

no sistema $F(2, 4, -3, 3)$.

Exercício 2.11

Sejam $x = 5/4$, $y = 3/8$. Represente $z = x + y$ no sistema $F(2, 3, -3, 2)$.

Exercício 2.12

Resolver a seguinte equação,

$$3x^2 - 2x - 1 = 0$$

e represente as soluções no sistema $F(10, 2, -2, 2)$.

Exercício 2.13

A derivada, $f'(x)$ de uma função $f(x)$ pode ser aproximada pela relação

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}$$



se $f(x) = 0,7e^{0,5x}$ e $h = 0,3$. Calcule

- o valor aproximado de $f'(2)$
- o erro relativo.

Exercício 2.14

Qual o menor valor de x de modo que,

$$\frac{e^x + \cos x - \sin x - 2}{x^3}$$

pode ser representado no sistema $F(2, 4, -5, 8)$?

Exercício 2.15

Represente o valor obtido por meio da seguinte expressão,

$$a + \sqrt{a^2 + b^3}$$

com $a = -123456$ e $b = 132$ no sistema $F(2, 4, -5, 6)$?



EXPANSÃO EM SÉRIES DE TAYLOR

A grande maioria neste curso já estudou séries infinitas (principalmente as séries de Taylor) em cálculo diferencial. Na Computação Numérica, iremos adquirir uma boa compreensão do uso prático deste tópico. Conseqüentemente, esta seção é particularmente importante em algumas aplicações e conceitos. Portanto, merece um estudo cuidadoso.

POLINÔMIO DE TAYLOR

Uma **série de potências** é uma expressão da forma,

$$\sum_{k=0}^{\infty} a_k x^k = a_0 + a_1 x + a_2 x^2 + a_3 x^3 + \dots$$

Mostraremos ao longo deste capítulo, como representar uma gama de funções por meio de séries de potências. O problema se resumirá a obtenção de uma aproximação para uma dada função. Ou seja,

$$y \approx f(x).$$

Para isso, considera-se uma função $f: I \subset \mathbb{R} \rightarrow \mathbb{R}$ derivável num ponto $x_0 \in I$, sendo I um intervalo aberto em \mathbb{R} e r uma reta tangente ao gráfico da função no ponto x_0 . Para valores de x numa vizinhança de x_0 , temos a aproximação linear $T: \mathbb{R} \rightarrow \mathbb{R}$ representando uma aproximação da função $f(x)$ por meio da reta r . Como r é a reta



que passa pelo ponto $(x_0, f(x_0))$, temos que sua equação pode ser obtida por meio da relação $f(x) - f(x_0) = \alpha(x - x_0)$, onde α representa o coeficiente angular ou a derivada de f no ponto x_0 . De modo que,

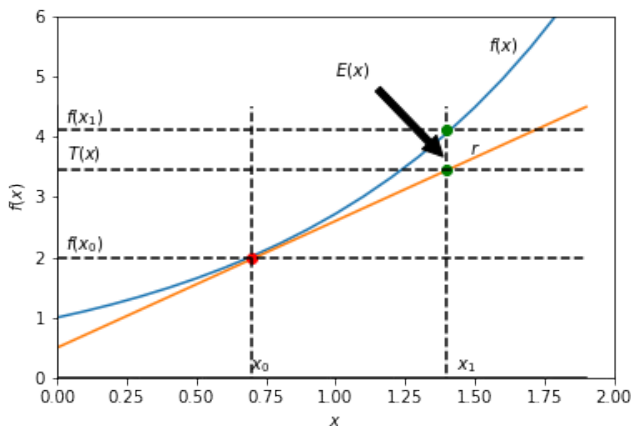
$$T(x) = f(x_0) + f'(x_0)(x - x_0).$$

Quando $x \rightarrow x_0$ temos, $T(x_0) = f(x_0) + f'(x_0)(x_0 - x_0) = f(x_0)$. Agora, considerando $x \in I$ um ponto na vizinhança de x_0 de maneira que $T(x) \neq f(x)$. De modo que tenhamos $E(x) = f(x) - T(x)$ o erro da aproximação da função $f(x)$ pela aproximação linear $T(x)$. De modo que,

$$\lim_{x \rightarrow x_0} E(x) = 0$$

e visualizado na seguinte figura.

Figura 3.1 – Gráfico representando o erro obtido via aproximação linear



De fato,

$$\lim_{x \rightarrow x_0} E(x) = \lim_{x \rightarrow x_0} (f(x) - T(x)) = f(x_0) - T(x_0) = 0,$$



pois,

$$f(x_0) = T(x_0). \text{ Como } T(x) = f(x_0) + f'(x_0)(x - x_0) \text{ e } E(x) = f(x) - T(x),$$

temos que

$$E(x) = f(x) - f(x_0) - f'(x_0)(x - x_0).$$

Para $x \neq x_0$,

$$\frac{E(x)}{(x - x_0)} = \frac{f(x) - f(x_0)}{x - x_0} - f'(x_0).$$

Além disso,

$$\lim_{x \rightarrow x_0} \frac{E(x)}{(x - x_0)} = \lim_{x \rightarrow x_0} \left(\frac{f(x) - f(x_0)}{x - x_0} - f'(x_0) \right) = f'(x_0) - f'(x_0) = 0.$$

Portanto, como o limite tende a zero, segue que $E(x) \rightarrow 0$ mais rapidamente que $(x - x_0)$. Além disso, a reta tangente é a única que possui essa propriedade. Dito de outra forma, se uma função $f(x)$ for derivável até primeira ordem num ponto $x_0 \in I$, podemos inferir a existência de um polinômio

$$p_1(x) = T(x) = f(x_0) + f'(x_0)(x - x_0).$$

Que constitui o **polinômio de Taylor de ordem 1**.

Em geral, precisamos de uma expressão que nos forneça o erro cometido num processo de aproximação, sem necessariamente ter que calcular o valor da função $f(x)$ num dado ponto. Para isso, temos:



Teorema 3.1

Seja $f: I \subset \mathbb{R} \rightarrow \mathbb{R}$ uma função derivável até segunda ordem no intervalo aberto I com x e $x_0 \in I$. Então, existe pelo menos um $x^- \in I$ tal que $x_0 < x^- < x$ e, $f''(x^-)$

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(\bar{x})}{2}(x - x_0)^2$$

Onde

$$E(x) = \frac{f''(\bar{x})}{2}(x - x_0)^2.$$

Demonstração.

De fato, sendo $E(x) = f(x) - p_1(x)$, temos

$$E(x) = f(x) - (f(x_0) + f'(x - x_0)) = f(x) - f(x_0) - f'(x - x_0).$$

Como $E(x_0) = 0$ e escolhendo $g(x) = (x - x_0)^2$ e $g(x_0) = 0$. Além disso,

$$\frac{E(x)}{g(x)} = \frac{E(x) - E(x_0)}{g(x) - g(x_0)}.$$

Então pelo Teorema do Valor Médio de Cauchy, existe um $x_0 < c < x$ tal que:

$$\frac{E(x)}{g(x)} = \frac{E(x) - E(x_0)}{g(x) - g(x_0)} = \frac{E'(c)}{g'(c)}.$$



Como, $E'(x) = f'(x) - f'(x_0)$ e $g'(x) = 2(x - x_0)$ segue que $E'(x_0) = 0$ e $g'(x_0) = 0$. Assim, obtemos a igualdade:

$$\frac{E(x)}{g(x)} = \frac{E'(c) - E'(x_0)}{g'(c) - g'(x_0)}$$

Aplicando novamente o Teorema do Valor Médio de Cauchy, existirá $x_0 < x^- < c$ tal que:

$$\frac{E'(c) - E'(x_0)}{g'(c) - g'(x_0)} = \frac{E''(\bar{x})}{g''(\bar{x})}$$

Consequentemente

$$\frac{E(x)}{g(x)} = \frac{E''(\bar{x})}{g''(\bar{x})}$$

sendo $E''(x) = f''(x)$ e $g''(x) = 2$ temos que $E''(x^-) = f''(x^-)$ e $g''(x^-) = 2$, portanto

$$\frac{E(x)}{g(x)} = \frac{f''(\bar{x})}{2} \implies E(x) = \frac{f''(\bar{x})}{2}(x - x_0)^2.$$

Exemplo 3.1

Seja $f: \mathbb{R}_+ \rightarrow \mathbb{R}$ definida como $f(x) = \sqrt{x}$. Utilize o polinômio de Taylor de ordem 1 para estimar o valor de $f(9, 02)$.

Devemos construir o polinômio de Taylor de ordem 1 da função $f(x) = \sqrt{x}$ considerando

uma vizinhança de $x_0 = 9$. Como $f(x_0) = f(9) = 3$ e $f'(x) = (2\sqrt{x})^{-1}$, temos

$$f'(x_0) = f'(9) = \frac{1}{2\sqrt{9}} = \frac{1}{6}.$$



De modo que o polinômio de Taylor de ordem 1 da função $f(x) = \sqrt[3]{x}$ em torno de uma vizinhança de $x_0 = 9$, pode ser escrito como:

$$p_1(x) = f(x_0) + f'(x_0)(x - x_0) = 3 + \frac{1}{6}(x - 9) = \frac{3}{2} + \frac{x}{6}.$$

$$p_1(9,02) = \frac{3}{2} + \frac{9,02}{6} = 3,0033.$$

Sabemos que $E_r(x) = |f(x) - p_1(x)|$. Portanto,

$$|E_r(9,02)| = \left| \sqrt[3]{9,02} - \left(\frac{3}{2} + \frac{9,02}{6} \right) \right| \approx 3,1483 \times 10^{-5}$$

Aplicando o resultado obtido no Teorema, temos que $f''(x) = -(4\sqrt[3]{x^3})^{-1}$. Portanto, para algum $x_0 < x^- < x$, temos

$$|E(9,02)| = \left| \frac{f''(\bar{x})}{2} (9,02 - 9)^2 \right| = \left| -\frac{1}{8\sqrt[3]{\bar{x}^3}} (0,02)^2 \right|.$$

Não podemos calcular x^- , mas sabemos que

$$\begin{cases} 9 < \bar{x} < 9,03 \\ 9^3 < \bar{x}^3 < 9,03^3 \\ \sqrt[3]{9^3} < \sqrt[3]{\bar{x}^3} < \sqrt[3]{9,03^3} \\ 3^3 < \sqrt[3]{\bar{x}^3} < \sqrt[3]{9,03^3} \\ \frac{1}{3^3} > \frac{1}{\sqrt[3]{\bar{x}^3}} > \frac{1}{\sqrt[3]{9,03^3}} \end{cases}$$

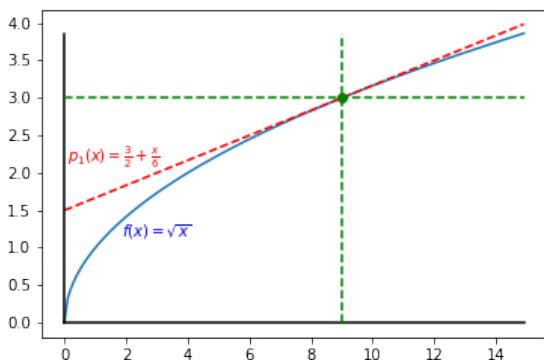
Da última desigualdade podemos concluir que

$$|E(9,02)| = \left| -\frac{1}{8\sqrt[3]{\bar{x}^3}} (0,02)^2 \right| < \left| -\frac{1}{8,3^3} (0,02)^2 \right| < 10^{-5}.$$



Assim, ao adotarmos $p_1(9, 02)$ estaremos cometendo um erro menor que 10^{-5} na aproximação.

Figura 3.2 – Gráfico aproximado de $f(x) = \sqrt{x}$ por um polinômio de Taylor de ordem 1



CASO GERAL

De maneira geral, supondo que a função f possui n derivadas no ponto x_0 . Existe então, um polinômio $p_n(x)$ de grau menor ou igual a n , tal que

$$f(x) = p_n(x) + \mathcal{O}((x - x_0)^n), \quad x \rightarrow x_0, \quad (3.2)$$

Por analogia, a partir da idéia apresentada na equação (3.2), podemos escrever o polinômio de grau n , por meio da seguinte série de potências

$$p_n(x) = a_0 + a_1(x - x_0) + a_2(x - x_0)^2 + \cdots + a_n(x - x_0)^n + \mathcal{O}((x - x_0)^n) \quad (3.3)$$

Considerando $\mathcal{O}((x - x_0)^n) \rightarrow 0$ e valorando o polinômio no ponto x_0 temos que



$$p_n(x_0) = a_0 + a_1(x_0 - x_0) + a_2(x_0 - x_0)^2 + \cdots + a_n(x_0 - x_0)^n \rightarrow a_0 = p_n(x_0) = f(x_0).$$

Derivando $p_n(x)$ com relação a x , obtemos:

$$p'_n(x) = a_1 + 2a_2(x - x_0) + \cdots + na_n(x - x_0)^{n-1}$$

valorando o polinômio no ponto x_0 temos que $a_1 = p'_n(x_0) = f'(x_0)$.

A derivada segunda de $p_n(x)$ com relação a x é,

$$p''_n(x) = 2.1.a_2 + \cdots + n(n-1)a_n(x - x_0)^{n-2} \rightarrow a_2 = \frac{p''_n(x_0)}{2!} = \frac{f''(x_0)}{2!}.$$

De maneira geral,

$$a_k = \frac{f^{(k)}(x_0)}{k!}, \quad k = 0, 1, \dots, n.$$

Ou seja,

$$\begin{cases} a_0 = p_n(x_0) = f(x_0) \\ a_1 = p'_n(x_0) = f'(x_0) \\ \quad \quad \quad \vdots \\ a_n = p_n^{(n)}(x_0) = f^{(n)}(x_0). \end{cases} \quad (3.4)$$

Assim, podemos escrever o polinômio $p_n(x)$ como,

$$p_n(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \cdots + \frac{f^{(k)}(x_0)}{k!}(x - x_0)^k + \cdots + \frac{f^{(n)}(x_0)}{n!}(x - x_0)^n. \quad (3.5)$$

que corresponde ao **polinômio de Taylor de grau n** . Cumprindo todas as relações apresentadas em (3.4). No entanto, dependendo do grau do polinômio, pode aparecer um **erro residual**. Ou seja,



$$E_r(x) = f(x) - p_n(x) = \mathcal{O}((x - x_0)^{n+1}) \approx \frac{f^{(n+1)}(\bar{x})}{(n+1)!} (x - x_0)^{n+1}, \quad x_0 < \bar{x} < x. \quad (3.6)$$

Obtido por meio do Teorema do Valor Médio de Cauchy.

Exemplo 3.2

Considere a função:

$$f(x) = 3x^5 - 2x^4 + 15x^3 + 13x^2 - 12x - 5.$$

Determine o polinômio de Taylor no ponto $x_0 = 2$.

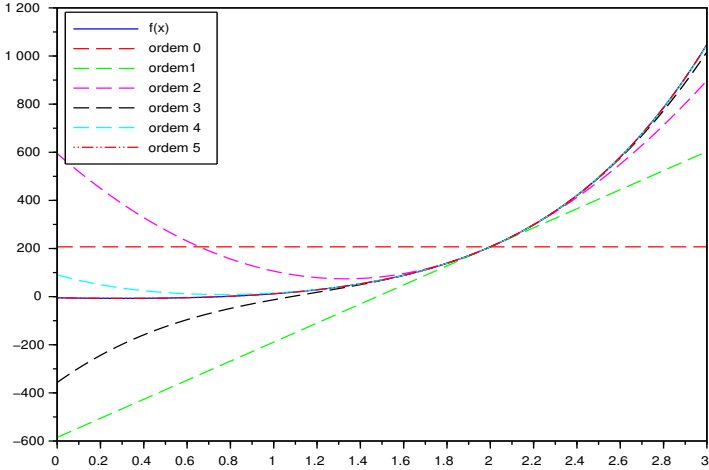
- ordem zero: $f(x) = 3x^5 - 2x^4 + 15x^3 + 13x^2 - 12x - 5 \rightarrow f(2) = 207$
- ordem um: $f'(x) = 15x^4 - 8x^3 + 45x^2 + 26x - 12 \rightarrow f'(2) = 396$
- ordem dois: $f''(x) = 60x^3 - 24x^2 + 90x + 26 \rightarrow f''(2) = 590$
- ordem três: $f'''(x) = 180x^2 - 48x + 90 \rightarrow f'''(2) = 714$
- ordem quatro: $f^{(4)}(x) = 360x - 48 \rightarrow f^{(4)}(2) = 672$
- ordem cinco: $f^{(5)}(x) = 360 \rightarrow f^{(5)}(2) = 360$
- ordem superior: $f^{(k)}(x) = 0 \rightarrow f^{(k)}(2) = 0$.

Portanto, para $k \geq 6$, temos:

$$\begin{aligned} f(x) &= f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots + \frac{f^{(k)}(x_0)}{k!}(x - x_0)^k \\ &= 207 + 396(x - 2) + \frac{590}{2!}(x - 2)^2 + \frac{714}{3!}(x - 2)^3 + \frac{672}{4!}(x - 2)^4 + \frac{360}{5!}(x - 2)^5 + 0 \\ &= 207 + 396(x - 2) + 295(x - 2)^2 + 119(x - 2)^3 + 28(x - 2)^4 + 3(x - 2)^5 \end{aligned}$$



Figura 3.3 – Aproximação utilizando série de Taylor da função $f(x) = 3x^5 - 2x^4 + 15x^3 + 13x^2 - 12x - 5$, com $0 \leq x \leq 3$, em torno do ponto $x_0 = 2$. Com os gráficos da ordem zero a cinco



No entanto, quando $n \rightarrow \infty$ na equação (3.2) e supondo que a função $f(x)$ definida num intervalo (a, b) é infinitamente diferenciável num ponto $x_0 \in (a, b)$. Então, quando $x \rightarrow x_0$.

$$f(x) = \sum_{k=0}^{\infty} \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k. \quad (3.7)$$

Que corresponde a **série de Taylor**. Podemos observar que o polinômio $p_n(x)$ definido anteriormente, representa a série de Taylor truncada em algum termo. Para simplificar a notação, iremos considerar o caso especial, em que $x_0 = 0$, de modo que



$$f(x) = \sum_{n=0}^{\infty} a_n x^n = a_0 + a_1 x + a_2 x^2 + \dots + a_k x^k + \dots$$

$$f'(x) = a_1 + 2a_2 x + 3a_3 x^2 + \dots + k a_k x^{(k-1)} + \dots$$

$$f''(x) = 2a_2 + 2.3 a_3 x + 2.3.4 a_4 x^2 + \dots + (k-1)k a_k x^{(k-2)} + \dots$$

$$\vdots \quad \ddots \quad \vdots$$

$$f^k(x) = 1.2 \dots k a_k + 2.3 \dots (k+1) a_{k+1} x + \dots + (n-k+1)(n-k+2) \dots k a_k x^{(n-k)} + \dots$$

Valorando a função e suas respectivas derivadas na origem. Ou seja, $x \rightarrow x_0 = 0$, obtemos:

$$\left\{ \begin{array}{l} f(0) = a_0 \\ f'(0) = a_1 \rightarrow a_1 = f'(0)/1! \\ f''(0) = 1.2 a_2 \rightarrow a_2 = f''(0)/2! \\ \vdots \\ f^{(k)}(0) = 1.2.3 \dots k a_k \rightarrow a_k = f^{(k)}(0)/k! \end{array} \right. \quad (3.8)$$

Portanto, a série de Taylor para $x_0 = 0$ pode ser escrita como,

$$f(x) = \sum_{k=0}^{\infty} \frac{f^{(k)}(0)}{k!} (x)^k. \quad (3.9)$$

Correspondendo a um caso particular, denominado de **série de Maclaurin**.

Exemplo 3.3

Obtenha a série de Taylor da função $f(x) = e^x$ em torno do ponto $x_0 = 0$.



$$\begin{aligned}
 f(x) &= e^x \rightarrow f(0) = 1 \\
 f'(x) &= e^x \rightarrow f'(0) = 1 \\
 &\vdots \\
 f^{(n)}(x) &= e^x \rightarrow f^{(n)}(0) = 1 \\
 &\vdots
 \end{aligned}$$

Assim,

$$e^x = 1 + x + \frac{x^2}{2!} + \cdots + \frac{x^n}{n!} + \cdots = \sum_{k=0}^{\infty} \frac{x^k}{k!}. \quad (3.10)$$

Se considerarmos um número finito de termos, podemos aproximar a função por:

$$e^x \approx \sum_{k=0}^n \frac{x^k}{k!} + \frac{e^{\xi}}{(n+1)!} x^{n+1}. \quad (3.11)$$

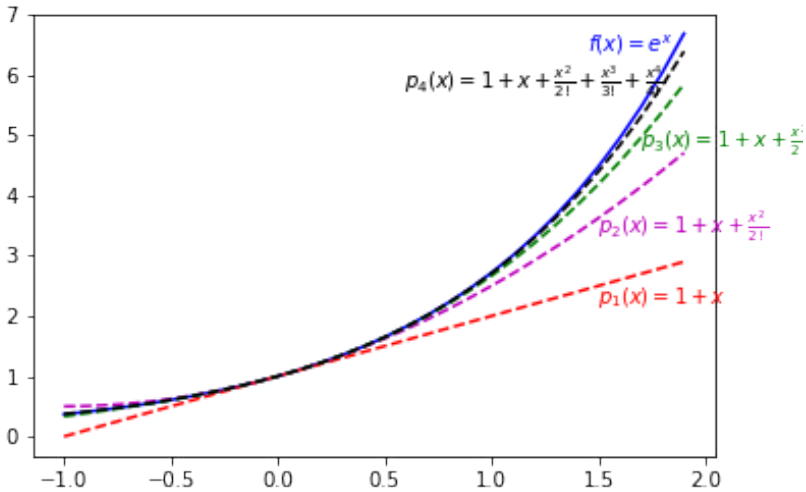
Considerando todos os valores de x pertencentes ao um intervalo simétrico em torno da origem, para $-\alpha \leq x \leq \alpha \Rightarrow |x| \leq \alpha$, $|\xi| \leq \alpha$, e $e^{\xi} \leq e^{\alpha}$. Assim, os termos restantes satisfazem a desigualdade:

$$\lim_{k \rightarrow \infty} \left| \frac{e^{\xi}}{(k+1)!} x^{k+1} \right| \leq \lim_{k \rightarrow \infty} \left| \frac{e^{\alpha}}{(k+1)!} \alpha^{k+1} \right| = 0.$$

Portanto, a série é convergente.



Figura 3.4 – Aproximação da função $f(x) = e^x$, em torno da origem. Com os gráficos dos polinômios de Taylor de ordem um a quatro



Exemplo 3.4

Calcule o valor de e .

Substituindo $x = 1$ na equação (3.10) temos,

$$e = 1 + 1 + \frac{1}{2!} + \frac{1}{3!} + \cdots + \frac{1}{k!} + \cdots .$$

Por se tratar de uma série infinita, iremos simplificá-la sem perder a acuracidade no cálculo. Para isso, iremos truncar a série a partir do nono termo, de modo que

$$e = 1 + 1 + \frac{1}{2!} + \frac{1}{3!} + \cdots + \frac{1}{9!} + E_9.$$

Onde E_9 representa um erro de truncamento da seguinte forma,



$$E_t = \frac{1}{10!} + \frac{1}{11!} + \frac{1}{12!} + \dots = \frac{1}{10!} \left(1 + \frac{1}{11} + \frac{1}{11 \cdot 12} + \frac{1}{11 \cdot 12 \cdot 13} + \dots \right)$$

observando os termos entre parênteses, vemos que os denominadores crescem numa ordem maior que 10. Ex: $11 > 10$, $11 \cdot 12 > 10^2$, $11 \cdot 12 \cdot 13 > 10^3$, etc. Portanto, sem perda de generalidade, iremos considerar,

$$E_t < \frac{1}{10!} \left(1 + \frac{1}{10} + \frac{1}{10^2} + \frac{1}{10^3} + \dots \right).$$

A soma dos termos entre parênteses, constituem a soma dos termos de uma progressão geométrica que converge para,

$$\sum_{k=0}^{\infty} \frac{1}{10^k} = 1 + \frac{1}{10} + \frac{1}{10^2} + \dots = \frac{1}{1 - \frac{1}{10}} = \frac{10}{9}.$$

De modo que,

$$E_t < \frac{1}{10!} \cdot \frac{10}{9} = \frac{1}{9 \cdot 9!} = \frac{1}{3265920} < 10^{-6}.$$

Exemplo 3.5

Obtenha a série de Taylor da função $f(x) = \cos(x)$ em torno do ponto $x_0 = 0$.

$$\begin{aligned} f(x) &= \cos(x) \rightarrow f(0) = 1 \\ f'(x) &= -\sin(x) \rightarrow f'(0) = 0 \\ f''(x) &= -\cos(x) \rightarrow f''(0) = -1 \\ f^{(3)}(x) &= \sin(x) \rightarrow f^{(3)}(0) = 0 \\ f^{(4)}(x) &= \cos(x) \rightarrow f^{(4)}(0) = 1 \\ f^{(5)}(x) &= -\sin(x) \rightarrow f^{(5)}(0) = 0 \\ &\vdots \end{aligned}$$



$$\cos(x) = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k}}{(2k)!}. \quad (3.12)$$

Exemplo 3.6

Obtenha a série de Taylor da função

$$f(x) = \frac{1}{1-x},$$

em torno do ponto $x_0 = 0$, com $|x| < 1$.

$$f(x) = (1-x)^{-1} \rightarrow f(0) = 1$$

$$f'(x) = (1-x)^{-2} \rightarrow f'(0) = 1$$

$$f''(x) = 2(1-x)^{-3} \rightarrow f''(0) = 2$$

$$f^{(3)}(x) = 6(1-x)^{-4} \rightarrow f^{(3)}(0) = 6$$

$$f^{(4)}(x) = 24(1-x)^{-5} \rightarrow f^{(4)}(0) = 24$$

$$f^{(5)}(x) = 120(1-x)^{-5} \rightarrow f^{(5)}(0) = 120$$

⋮

$$\frac{1}{1-x} = 1 + x + 2\frac{x^2}{2!} + 6\frac{x^3}{3!} + 24\frac{x^4}{4!} + 120\frac{x^5}{5!} + \dots = 1 + x + x^2 + x^3 + x^4 + \dots = \sum_{k=0}^{\infty} x^k. \quad (3.13)$$

R Algumas funções possuem uma fórmula fechada da expansão da série de Taylor, tais como: exponencial, seno, cosseno, seno hiperbólico, cosseno hiperbólico, arco tangente, etc.

Podemos concluir dizendo que a série de Taylor (MacLaurin) é uma decorrência natural do problema de aproximação de uma função $f(x)$ por meio de polinômios. Essa aproximação se baseia na obtenção



de um bom ajuste da função em torno de uma vizinhança de um dado ponto $x = x_0$. A série em geral converge em um intervalo contendo x_0 , porém qualquer uma das somas parciais é um polinômio que estará suficiente próximo de $f(x)$ numa vizinhança bem restrita de x_0 . Ou seja, a série de Taylor é um aproximador local. Aproximadores globais serão estudados em capítulos posteriores.

EXERCÍCIOS RESOLVIDOS

Utilizando a série de Maclaurin calcule o valor da função $f(x) = \cos(x)$, no ponto $x_1 = \pi/4$ a partir do ponto $x_0 = \pi/12$. Com erro relativo $\leq 10^{-3}$.

$$h = x_1 - x_0 = \frac{\pi}{4} - \frac{\pi}{12} = \frac{\pi}{6}$$

$$f(x+h) = f(x_0) + f'(x_0)h + \dots + \frac{f^{(k)}(x_0)}{k!}h^k + \dots$$

- ordem zero: $f(x) = \cos(x) \rightarrow f(\pi/12) = 0,9659$
- ordem um: $f'(x) = -\sin(x) \rightarrow f'(\pi/12) = -0,2588$.

Portanto a aproximação de ordem um é dada por,

$$Ord_1(x) \approx f(x_0) + f'(x_0)h = 0,9659 - 0,2588 \cdot \frac{\pi}{6} = 0,8304$$

- ordem dois: $f''(x) = -\cos(x) \rightarrow f''(\pi/12) = -0,9659$

a aproximação de ordem dois é dada por,

$$Ord_2(x) \approx Ord_1(x) - \frac{0,9659}{2} \cdot \left(\frac{\pi}{6}\right)^2 = 0,6980$$

no caso como estamos utilizando uma sequência numérica, o erro relativo é dado por,



$$E_r = \frac{|f(x_{i+1}) - f(x_i)|}{|f(x_{i+1})|}$$

Ou seja,

$$E_r = \frac{|0,6980 - 0,8304|}{|0,6980|} = 0,1897 > 10^{-3}$$

- ordem três: $f'''(x) = \text{sen}(x) \rightarrow f'''(\pi/12) = 0,2588$.

Portanto,

$$\text{Ord}_3(x) \approx \text{Ord}_2 + \frac{0,2588}{6} \cdot \left(\frac{\pi}{6}\right)^3 = 0,7042$$

$$E_r = \frac{|0,7042 - 0,6980|}{|0,7042|} = 0,0088 > 10^{-3}$$

como o erro está diminuindo, significa que a sequência está convergindo para a solução desejada.

- ordem quatro: $f^{(4)}(x) = -\cos(x) \rightarrow f^{(4)}(\pi/12) = -0,9659$.

Assim,

$$\text{Ord}_4(x) \approx \text{Ord}_3 - \frac{0,9659}{24} \cdot \left(\frac{\pi}{6}\right)^4 = 0,7012$$

$$E_r = \frac{|0,7012 - 0,7042|}{|0,7012|} = 0,0043 > 10^{-3}$$

- ordem cinco: $f^{(5)}(x) = \text{sen}(x) \rightarrow f^{(5)}(\pi/12) = 0,2588$.

Ou seja,

$$\text{Ord}_5(x) \approx \text{Ord}_4 + \frac{0,2588}{120} \cdot \left(\frac{\pi}{6}\right)^5 = 0,7013$$

$$E_r = \frac{|0,7013 - 0,7012|}{|0,7013|} = 1.2104 \cdot 10^{-4} < 10^{-3}$$



Assim, podemos concluir que para a tolerância de erro solicitada, a solução é dada por:

$$\cos\left(\frac{\pi}{4}\right) \approx \text{Ord}_5(x) = 0,7013.$$

Uma maneira mais rápida de resolução do exercício, é fazer uso da expansão obtida em (??):

$$\cos(x) = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k}}{(2k)!}$$

- ordem zero:

$$\cos\left(\frac{\pi}{4}\right) \approx \frac{\left(\frac{\pi}{4}\right)^0}{0!} = 1$$

- ordem um:

$$\cos\left(\frac{\pi}{4}\right) \approx \frac{\left(\frac{\pi}{4}\right)^0}{0!} - \frac{\left(\frac{\pi}{4}\right)^2}{2!} = 0,6916$$

- ordem dois:

$$\cos\left(\frac{\pi}{4}\right) \approx \frac{\left(\frac{\pi}{4}\right)^0}{0!} - \frac{\left(\frac{\pi}{4}\right)^2}{2!} + \frac{\left(\frac{\pi}{4}\right)^4}{4!} = 0,7074$$

- ordem três:

$$\cos\left(\frac{\pi}{4}\right) \approx \frac{\left(\frac{\pi}{4}\right)^0}{0!} - \frac{\left(\frac{\pi}{4}\right)^2}{2!} + \frac{\left(\frac{\pi}{4}\right)^4}{4!} - \frac{\left(\frac{\pi}{4}\right)^6}{6!} = 0,7078$$

$$E_r = \frac{|0,7078 - 0,7074|}{|0,7078|} = 4.6060 \cdot 10^{-4} < 10^{-3}$$



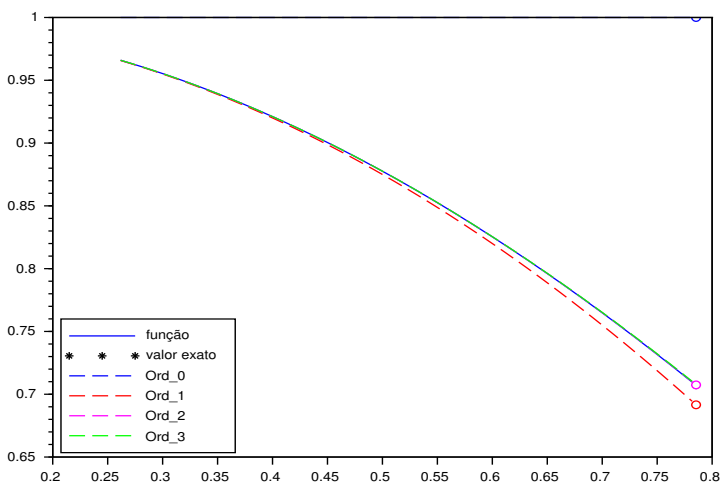
Assim, podemos concluir que para a tolerância de erro solicitada, a solução é dada por:

$$\cos\left(\frac{\pi}{4}\right) \approx \text{Ord}_3(x) = 0,7078.$$

Ou seja, desta forma não é necessário conhecer o intervalo Δx , basta conhecermos o ponto em que desejamos calcular o valor da função.

A seguinte figura, apresenta a convergência da série em direção ao valor desejado.

Figura 3.5 – Aproximação utilizando série de Taylor da função $f(x) = \cos(x)$, com $0,2 \leq x \leq 0,8$, para determinar $\cos(\pi/4)$. Com as aproximações da ordem zero a três



Exercício 3.2

A partir das funções $\text{sen}(x)$ e $\text{cos}(x)$ definidas por:



$$\operatorname{sen}(x) = \frac{e^{ix} - e^{-ix}}{2i}, \quad \operatorname{cos}(x) = \frac{e^{ix} + e^{-ix}}{2}.$$

Mostre que as séries de Taylor (MacLaurin) das mesmas são dadas por,

$$\operatorname{sen}(x) = \sum_{k=0}^{\infty} (-1)^k \frac{1}{(2k+1)!} x^{2k+1}, \quad \operatorname{cos}(x) = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k}}{(2k)!}$$

Sabemos que,

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots$$

substituindo x por ix , obtemos:

$$e^{ix} = \sum_{k=0}^{\infty} \frac{(ix)^k}{(ix)!} = 1 + ix - \frac{x^2}{2!} - i\frac{x^3}{3!} + \frac{x^4}{4!} + i\frac{x^5}{5!} - \frac{x^6}{6!} - i\frac{x^7}{7!} + \dots$$

$$e^{ix} = \left(1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots\right) + i \left(x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots\right)$$

$$e^{ix} = \left(\sum_{k=0}^{\infty} (-1)^k \frac{x^{2k}}{(2k)!}\right) + i \left(\sum_{k=0}^{\infty} (-1)^k \frac{x^{2k+1}}{(2k+1)!}\right) \quad (3.14)$$

e, temos que

$$e^{-x} = \sum_{k=0}^{\infty} (-1)^k \frac{x^k}{k!} = 1 - x + \frac{x^2}{2!} - \frac{x^3}{3!} + \frac{x^4}{4!} - \frac{x^5}{5!} + \dots$$



substituindo x por $-ix$, obtemos:

$$e^{-ix} = \sum_{k=0}^{\infty} (-1)^k \frac{(ix)^k}{(ix)^k} = 1 - ix - \frac{x^2}{2!} + i\frac{x^3}{3!} + \frac{x^4}{4!} - i\frac{x^5}{5!} + \dots$$

$$e^{-ix} = \left(1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots\right) - i \left(x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots\right)$$

$$e^{-ix} = \left(\sum_{k=0}^{\infty} (-1)^k \frac{x^{2k}}{(2k)!}\right) - i \left(\sum_{k=0}^{\infty} (-1)^k \frac{x^{2k+1}}{(2k+1)!}\right) \quad (3.15)$$

somando as equações (3.14) e (3.15), temos que

$$e^{ix} + e^{-ix} = 2 \left(\sum_{k=0}^{\infty} (-1)^k \frac{x^{2k}}{(2k)!}\right) \rightarrow \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k}}{(2k)!} = \frac{e^{ix} + e^{-ix}}{2}$$

e subtraindo-as, obteremos

$$e^{ix} - e^{-ix} = 2i \left(\sum_{k=0}^{\infty} (-1)^k \frac{x^{2k+1}}{(2k+1)!}\right) \rightarrow \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k+1}}{(2k+1)!} = \frac{e^{ix} - e^{-ix}}{2i}$$

Exercício 3.3

Para datar rochas ou artefatos com mais de 50.000 anos, é preciso usar outros elementos radioativos. A equação seguinte é válida para qualquer isótopo radioativo:

$$\frac{S(t) - S(0)}{R(t)} + 1 = e^{(ln2)t/\lambda}$$



onde $S(t)$ representa o número de átomos do produto estável resultante do decaimento radioativo, $S(0)$ é o número inicial no produto estável, $R(t)$ o número de átomos do isótopo radioativo no instante t e λ é a meia vida do isótopo radioativo. Determine o valor aproximado do tempo usando um polinômio de Taylor de ordem 2 para a exponencial em torno da origem.

$$e^{(ln2)t/\lambda} \approx 1 + \frac{(ln2)t}{\lambda} + \frac{((ln2)t/\lambda)^2}{2!}$$

$$\frac{S(t) - S(0)}{R(t)} \approx \frac{(ln2)t}{\lambda} + \frac{((ln2)t/\lambda)^2}{2!}$$

$$\frac{(ln2)^2}{2!\lambda^2}t^2 + \frac{(ln2)}{\lambda}t - \frac{S(t) - S(0)}{R(t)} \approx 0$$

que constitui uma equação de segunda grau na variável t , cuja solução para $t > 0$ é dada por,

$$t \approx \frac{\lambda}{ln2} \left(\sqrt{1 + 2 \left(\frac{S(t) - S(0)}{R(t)} \right)} - 1 \right).$$

Exercício 3.4

Obtenha a expansão em série de Taylor da função

$$f(x) = x^6 e^{2x^3}$$



em torno da origem.

$$x^6 e^{2x^3} = x^6 \sum_{k=0}^{\infty} \frac{(2x^3)^k}{k!} = \sum_{k=0}^{\infty} x^6 \frac{2^k x^{3k}}{k!} = \sum_{k=0}^{\infty} \frac{2^k x^{3k+6}}{k!}$$

EXERCÍCIOS PROPOSTOS

Exercício 3.5

Obtenha a série de Taylor da função $f(x) = \text{sen}(x)$ em torno do ponto $x_0 = 0$.

Exercício 3.6

Calcule $f(x) = \text{sen}(x)$ para $x = 10^\circ$ com $E_t < 10^{-5}$.

Exercício 3.7

Obtenha a série de Taylor da função $f(\theta) = e^{i\theta} \in \mathbb{C}$ em torno do ponto $\theta_0 = 0$ em radianos e, sabendo que $i^2 = -1$ é a unidade imaginária.



Exercício 3.8

Sabemos do cálculo diferencial que a função $\tan : (-\pi/2, \pi/2) \rightarrow \mathbb{R}$ é uma bijeção em C^∞ com derivada positiva, e sua inversa $\tan^{-1} : \mathbb{R} \rightarrow (-\pi/2, \pi/2)$ possui derivada igual a $1/(1+x^2), \forall x \in \mathbb{R}$. A partir das informações apresentadas obtenha a série de Taylor da função $\tan^{-1}(x)$.

Exercício 3.9

Mostre que a relação,

$$\frac{\pi}{4} = \tan^{-1}\left(\frac{1}{2}\right) + \tan^{-1}\left(\frac{1}{3}\right)$$

é verdadeira.

Exercício 3.10

Obtenha a série de Taylor da função $f(x) = (1+x)^\gamma$, onde $\gamma \in \mathbb{R}$.

Exercício 3.11

Para Calcular as coordenadas espaciais de um planeta, é necessário calcular a função



$$f(x) + 1 = x - 0,5 \operatorname{sen} x$$

tome como base o ponto $x_i = \pi$ no intervalo $[0; \pi]$. Determine a expansão em série de Taylor de ordem mais alta que possa ser representada no sistema $F(2, 5, -4, 4)$ com erro relativo máximo de 5%.

Exercício 3.12

Encontre o polinômio de Taylor de ordem 2 que representa a função f na origem do sistema de coordenadas, onde a função $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ é dada por:

$$f(x, y) = x \operatorname{sen}(y), \quad (x, y) \in \mathbb{R}^2.$$

Exercício 3.13

Utilizando a expansão em série de Taylor, demonstre que:

$$f(x, y) = x \operatorname{sen}(y), \quad (x, y) \in \mathbb{R}^2.$$

Exercício 3.14

Encontre via série de Taylor, um polinômio com coeficientes reais e grau ≤ 5 , de modo que



$$\int_{-\pi}^{\pi} |\text{sen}(x) - p_n(x)|^2 dx \rightarrow 0.$$

Exercício 3.15

Utilizando expansão em série de Taylor, calcule:

$$I = \int_0^1 e^{-t^2} dt, \text{ com } E_a \leq 10^{-2}$$

Exercício 3.16

Utilizando expansão em série de Taylor, calcule:

$$I = \int_{0,2}^{0,4} \frac{\ln(1+x)}{x} dx, \text{ com } E_a \leq 10^{-3}$$

Exercício 3.17

Utilizando a expansão em série de Taylor, calcule

$$\int_0^1 \frac{\text{sen}(\pi x)}{\pi x} dx$$

Exercício 3.18

As funções $\text{senh}(x)$ e $\text{cosh}(x)$ são definidas por:



$$\operatorname{senh}(x) = \frac{e^x - e^{-x}}{2}, \quad \operatorname{cosh}(x) = \frac{e^x + e^{-x}}{2}.$$

Mostre que as séries de Taylor das mesmas são dadas por,

$$\operatorname{senh}(x) = \sum_{k=0}^{\infty} \frac{1}{(2k+1)!} x^{2k+1}, \quad \operatorname{cosh}(x) = \sum_{k=0}^{\infty} \frac{1}{(2k)!} x^{2k}$$

Exercício 3.19

As funções $\operatorname{sen}(x)$ e $\operatorname{cos}(x)$ são definidas por:

$$\operatorname{sen}(x) = \frac{e^{ix} - e^{-ix}}{2i}, \quad \operatorname{cos}(x) = \frac{e^{ix} + e^{-ix}}{2}.$$

Mostre que as séries de Taylor das mesmas são dadas por,

$$\operatorname{sen}(x) = \sum_{k=0}^{\infty} (-1)^k \frac{1}{(2k+1)!} x^{2k+1}, \quad \operatorname{cos}(x) = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k}}{(2k)!}$$

Exercício 3.20

Encontre uma aproximação via séries de MacLaurin para a seguinte função,

∞



$$f(x) = e^x \cos x + \ln \left(x + \frac{1}{2} \right)$$

com uma tolerância de 10^{-3} .

Exercício 3.21

Calcule o seguinte limite,

$$\lim_{x \rightarrow 0} \left(\frac{1}{\text{sen}(x)} - \frac{1}{x} \right)$$

por meio das séries de Maclaurin.



SOLUÇÃO DE EQUAÇÕES NÃO LINEARES

Este capítulo é devotado aos problemas de determinar raízes (ou zeros) de equações, por meio do estudo numérico de alguns métodos iterativos na resolução de tais problemas (ou equações não lineares). Em Ciência e Engenharia muitas vezes necessitamos resolver problemas por meio de encontrar a raiz de equações algébricas, transcendentais ou não lineares e que geralmente nenhuma informação sobre a raiz está disponível. Em geral, para encontrar tal valor numérico, desejamos encontrar alguns valores aproximados que atendem a requisitos sem necessariamente comprometer a solução final. Eis um exemplo de um problema que ocorre frequentemente em trabalhos científicos e/ou voltados a resolver alguns casos na Engenharia.

Por exemplo, A área S da superfície lateral de um cone é dada por:

$$S = \pi r \sqrt{r^2 + h^2}$$

onde r representa o raio da base e h a altura. Para determinar o raio do cone, uma forma é resumir o problema a encontrar um valor de $r > 0$, que satisfaça a equação

$$f(r) = 0. \tag{4.1}$$

No caso, $f(r) = S - \pi r \sqrt{r^2 + h^2} = 0$.



Em resumo, o principal objetivo é escrever o problema como uma função $f: \mathbb{R} \rightarrow \mathbb{R}$, de modo a encontrar os valores de x , que satisfaçam a relação $f(x) = 0$ (nos reais).

Antes de apresentar alguns métodos numéricos de resolução para tais problemas, eis um teorema de extrema importância a ser aplicado em alguns casos para isolar a raiz num determinado intervalo fechado em \mathbb{R} .

Teorema 4.1

Teorema do Valor Intermediário ou Teorema de Bolzano.
Seja $f: [a, b] \rightarrow \mathbb{R}$ uma função contínua no intervalo, de modo que

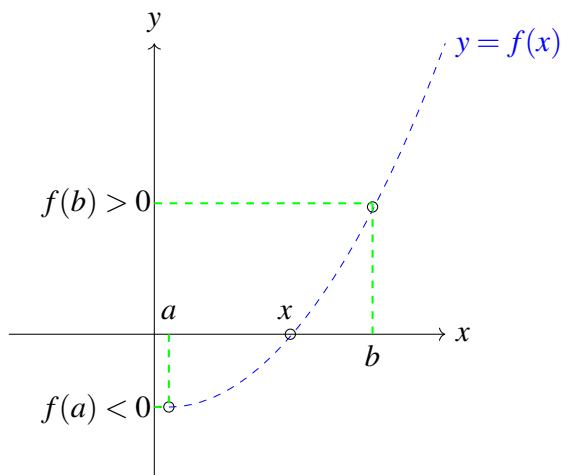
$$f(a) \cdot f(b) < 0 \tag{4.2}$$

então

$$\exists x \in [a, b] \Rightarrow f(x) = 0. \tag{4.3}$$



Figura 4.1 – Exemplo de aplicação do Teorema de Bolzano



O teorema do valor intermediário ou teorema de Bolzano (por vezes chamado teorema de Bolzano-Cauchy) estabelece que, se uma função f é contínua em um intervalo fechado $[a, b]$ e assume valores com sinais opostos em dois pontos distintos dentro desse intervalo, então existe pelo menos um ponto x contido no intervalo de modo que $f(x) = 0$. O Teorema de Bolzano só pode ser aplicado em funções contínuas num intervalo. Se a função não for contínua, o teorema não se aplica. Atenção, o teorema não explicita que existe um único ponto no intervalo de modo que sua imagem seja nula, apenas afirma que existe pelo menos um ponto, mas este pode não ser o único.

Exemplo 4.1

Isole as raízes da função $f(x) = x^3 - 6x + 2$. Claramente, pelo gráfico.



$$f(x) = \begin{cases} 0, & \text{para algum } x \in [-3, -2] \rightarrow f(-3) = -7 \text{ e } f(-2) = 6 \\ 0, & \text{para algum } x \in [0, 1] \rightarrow f(0) = 2 \text{ e } f(1) = -3 \\ 0, & \text{para algum } x \in [2, 3] \rightarrow f(2) = -2 \text{ e } f(3) = 11 \end{cases} \quad (4.4)$$

Figura 4.2 - Gráfico da função $f(x) = x^3 - 6x + 2$, mostrando a localização de suas 03 (três) raízes



Uma vez que vimos como isolar uma raiz, o próximo passo consiste em refinar a busca no intervalo, de modo a obtermos a solução (raiz) desejada. Ou seja, encontrar soluções para a equação não linear univariável $f(x) = 0$ contínua em $f: [a; b] \subset \mathbb{R} \rightarrow \mathbb{R}$. Para isso, iremos apresentar alguns métodos de resolução de tais equações.

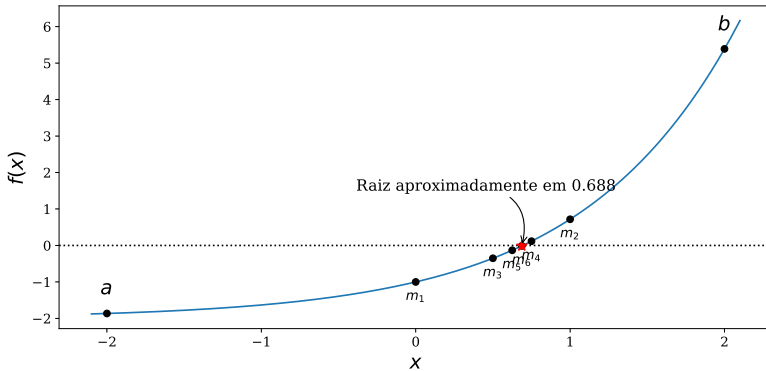
MÉTODO DA BISSECÇÃO

O primeiro método apresentado será o da bissecção, que na verdade é uma ideia maravilhosamente simples, baseada em pouco mais do que a continuidade da função e extremamente eficiente, porém com um custo computacional elevado. O método consiste em aproximações sucessivas por meio da divisão do intervalo $[a,$



$b]$ em subintervalos que contenham a raiz desejada, por meio da utilização do Teorema de Bolzano.

Figura 4.3 – Visualização gráfica do método da bissecção para encontrar a raiz da função $f(x) = e^x - 2$ no intervalo de $[-2, 2]$



Uma estimativa do número de passos necessários ao cálculo da raiz pode ser obtida por meio da seguinte expressão,

$$n \geq \frac{1}{\log(2)} \left(\log \left(\frac{b-a}{tol} \right) \right)$$

e a condição de parada definida por,

$$E_r = \frac{|x_{i+1} - x_i|}{|x_{i+1}|} \leq tol.$$

Eis, de maneira intuitiva o algoritmo do método:



Algorithm 1 Método da Bissecção

Passo 1: entradas (a, b, tol) de modo que $f(a) \cdot f(b) < 0$.

Passo 2 : Uma estimativa da raiz é determinada por:

$$x_r = \frac{a+b}{2}$$

if $f(a) \cdot f(x_r) < 0$ **then**

$b = x_r$

 volte ao passo 2

else

$a = x_r$

end if

if $E_r \leq tol$ **then**

x_r é solução

end if

onde tol representa a tolerância definida pelo usuário.

Exemplo 4.2

Utilizando o algoritmo da bissecção. Quantos passos são necessários para calcular a raiz de uma função f numa máquina de mantissa $p = 32$ com uma tolerância de $tol = 2^{-101002}$? Considere $a = 10000_2$ e $b = 10001_2$.

$$\begin{cases} 2^{-101002} = 2^{-2010} \\ a = 10000_2 = 16_{10} \Rightarrow n \geq \frac{1}{\log(2)} \left(\log \left(\frac{17-16}{2^{-20}} \right) \right) = \frac{1}{\log(2)} (20 \cdot \log(2)) = 20 \rightarrow n \geq 20. \\ b = 10001_2 = 17_{10} \end{cases}$$

Exercício 4.1

Usando o método da bissecção, resolva a equação



$$x^2 = -\ln(x)$$

com uma tolerância menor ou igual a 0,02.

Escolhendo

$$a = 0,5 \rightarrow f(a) = -0,4431 < 0$$

$$b = 1 \rightarrow f(b) = 1 > 0$$

$$x_r \in [a, b] = [0,5; 1] \Rightarrow x_r = 0,75.$$

Tabela 4.1 – Solução do exercício 4.1.

n	a	b	x_r	$f(x_r)$	erro
0	0,5	1,00	0,75	0,2748	—
1	0,5	0,75	0,625	-0,0794	0,2000
2	0,625	0,75	0,6875	0,0980	0,0909
3	0,625	0,6875	0,6562	0,0093	0,0476
4	0,625	0,6562	0,6406	-0,0035	0,0244
5	0,6406	0,6562	0,6484	-0,0128	0,0120

Exercício 4.2

Num determinado circuito elétrico, a tensão V e a corrente I estão relacionadas por meio das seguintes equações:

$$\begin{cases} I = a(e^{bV} - 1) \\ c = dI + V \end{cases}$$

onde a ; b ; c e d são constantes. Considerando $a = 2$; $b = 2$; $c = 12$ e $d = 7$. Determine a tensão no circuito.



Exercício 4.3

O empréstimo de uma quantia D em reais, em um certo banco. Cobra uma taxa de juros anual de $r = 5\%$, por um período de n anos, sendo o pagamento mensal, M , inserido na seguinte equação,

$$D = \frac{12M}{r} \left[1 - \left(1 + \frac{r}{12} \right)^{-12n} \right]$$

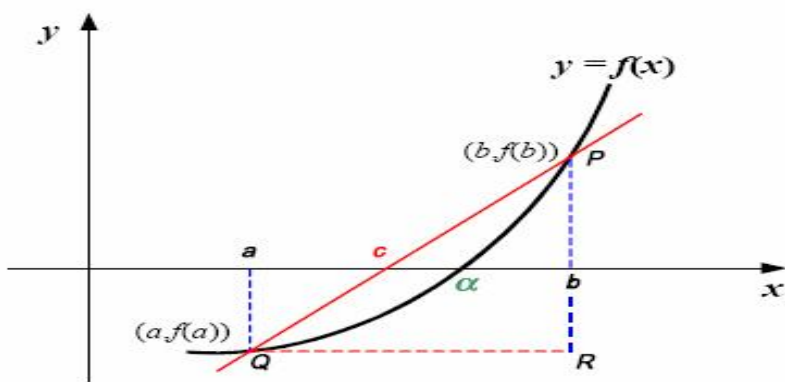
Um investidor necessita de um empréstimo de R\$160.000,00 para comprar um imóvel com uma prestação de R\$650,00 mensais. Supondo uma hipoteca de 35 anos, utilizando o método da bissecção determine a taxa de juros que o investidor irá pagar no período.

REGULA FALSI OU MÉTODO DA FALSA POSIÇÃO

O método da falsa posição é uma alternativa ao método da bissecção. Sabendo através do teorema de Bolzano que a raiz encontra-se no intervalo determinado. Podemos interligar os pontos $f(a)$ a $f(b)$ via uma reta. A intersecção dessa reta com o eixo x representa uma estimativa melhorada da raiz (conforme pode ser visto na figura (4.4)). O fato de substituirmos a curva por uma reta dar uma “falsa posição” da raiz. Eis a possível origem do nome do método. Utilizando semelhança de triângulos na citada figura,



Figura 4.4 – Exemplo método da falsa posição



podemos inferir uma equação para o cálculo do ponto de intersecção da reta com o eixo x . Ou seja,

$$\frac{f(a)}{c-a} = \frac{f(b)}{c-b} \rightarrow f(a)(c-b) = f(b)(c-a) \rightarrow c = \frac{bf(a) - af(b)}{f(a) - f(b)}$$

se,

$$\begin{aligned} c &= \frac{bf(a)}{f(a) - f(b)} - \frac{af(b)}{f(a) - f(b)} \\ c &= b + \frac{bf(a)}{f(a) - f(b)} - b - \frac{af(b)}{f(a) - f(b)} \\ c &= b + \frac{bf(a)}{f(a) - f(b)} - b - \frac{af(b)}{f(a) - f(b)} \\ c - b &= \frac{bf(a)}{f(a) - f(b)} - b \frac{f(a) - f(b)}{f(a) - f(b)} - \frac{af(b)}{f(a) - f(b)} \\ c - b &= \frac{bf(a) - bf(a) + bf(b) - af(a)}{f(a) - f(b)} \\ c &= b - \frac{f(b)(a-b)}{f(a) - f(b)} \end{aligned}$$



representada iterativamente por,

$$x_{i+1} = x_i - \frac{f(x_i)(x_{i-1} - x_i)}{f(x_{i-1}) - f(x_i)} \quad (4.5)$$

Embora o método da falsa posição seja preferível entre os métodos intervalares, existem algumas funções em que o método da bissecção apresenta melhores resultados. Por exemplo, $f(x) = x^{10} - 1$.

CONDIÇÃO DE CONVERGÊNCIA

Seja $\{x_k\}$ uma sequência de valores de x_r obtidos pelo método proposto e seja x^* a raiz exata de uma dada equação. Então o método é dito convergente se

$$\lim_{k \rightarrow \infty} |x_k - x^*| = 0.$$

ORDEM DE CONVERGÊNCIA

Seja E_n o erro obtido na n -ésima iteração com valor x_n , então

$$x_n = x^* + E_n, \quad x_{n+1} = x^* + E_{n+1},$$

que substituindo na equação (4.5), temos:

$$E_{n+1} = \frac{E_{n-1}f(x^* + E_n) - E_n f(x^* + E_{n-1})}{f(x^* + E_n) - f(x^* + E_{n-1})}$$

$$E_{n+1} = \frac{E_{n-1} [f(x^*) + E_n f'(x^*) + \dots] - E_n [f(x^*) + E_{n-1} f'(x^*) + \dots]}{[f(x^*) + E_n f'(x^*) + \dots] - [f(x^*) + E_{n-1} f'(x^*) + \dots]}$$

$$E_{n+1} = \frac{1}{2} E_{n-1} E_n \frac{f''(x^*)}{f'(x^*)} + \mathcal{O}(E_n^2), \quad f(x^*) = 0.$$



Determinando um número k tal que $E_{n+1} = ME^k$. E, portanto, $E_n = ME^k$.

MÉTODO DE NEWTON

Seja $f(x)$ uma função contínua no intervalo $[a; b]$ e x sua única raiz no intervalo; as derivadas $f'(x)$, ($f'(x) \neq 0$) e $f''(x)$ devem ser contínuas no intervalo. Encontra-se uma aproximação x_i para a raiz x feita pela expansão em série de Taylor para $f(x) = 0$:

$$f(x) \approx f(x_i) + f'(x_i)(x - x_i)$$

$$f(x_{i+1}) = f(x_i) + f'(x_i)(x_{i+1} - x_i) = 0$$

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}, i = 0, 1, \dots$$

onde: $x_{i+1} \approx x$.

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} \quad (4.6)$$

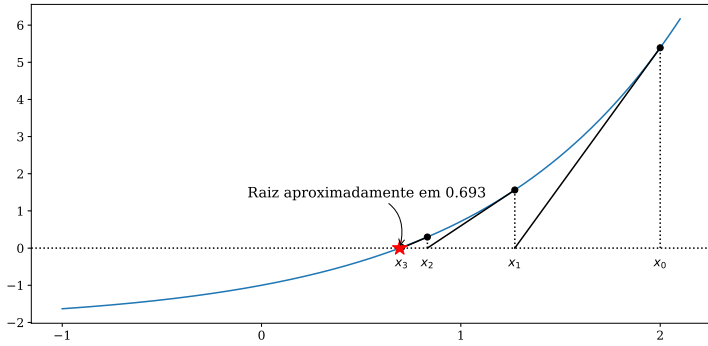
que é geralmente denominada **fórmula de Newton**.

É condição suficiente para a convergência do método de Newton que $f'(x)$ e $f''(x)$ sejam não nulas e preservem o sinal no intervalo $[a; b]$ e x_0 seja tal que:

$$f(x_0) \cdot f''(x_0) > 0$$



Figura 4.5 – Visualização gráfica do método de Newton para encontrar a raiz da função $f(x) = e^x - 2$ no intervalo de $[-1, 2]$



Algorithm 2 Método de Newton

Passo 1: Defina tol e escolha a condição inicial x_0 tal que,

$$f(x_0) \cdot f''(x_0) > 0.$$

Passo 2: para $i = 0, 1, 2, \dots, n$ faça:

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$$

if $E_r \leq tol$ **then**

x_{i+1} é a solução.

else

 volte ao passo 2.

end if



R Evita-se utilizar o método de Newton na resolução de equações cujo gráfico apresente pouca inclinação próxima a pontos de intersecção com o eixo x . Pois, assim

$$0 \rightarrow \frac{f(x)}{f'(x)} \rightarrow \infty.$$

Exemplo 4.3

Use o método de Newton para fazer uma estimativa da raiz de $f(x) = e^{-x} - x$, com $tol \leq 10^{-2}$.

Tabela 4.2 – Solução do exemplo 4.3

i	x_i	$f(x_i)$	$f'(x_i)$	x_{i+1}	erro
0	1	-0,6321	-1,3678	0,5378	—
1	0,5378	0,0461	-1,5839	0,567	0,0513
2	0,567	0,0002	-1,5672	0,5671	0,00027

Exemplo 4.4

Utilizando o método de Newton, calcule uma raiz da equação,

$$f(x) = e^x - 1,5 - tg^{-1}(x)$$

com $tol \leq 10^{-2}$

$$\begin{cases} f(x) = e^x - 1,5 - tg^{-1}(x) \\ f'(x) = e^x - (1+x^2)^{-1} \end{cases}$$



Tabela 4.3 – Solução do exemplo 4.4

i	x_i	$f(x_i)$	$f'(x_i)$	x_{i+1}	erro
0	1	0,4328	2,2183	0,8048	—
1	0,8048	0,0586	1,6295	0,7688	0,0468
2	0,7688	0,0018	1,5287	0,7676	0,0015

Exercício 4.4

Calcular $\sqrt{\alpha}$ com $\alpha > 0$, utilizando o método de Newton.

Fazendo,

$$x = \sqrt{\alpha} \rightarrow x^2 = \alpha \rightarrow \begin{cases} f(x) = x^2 - \alpha \\ f'(x) = 2x \end{cases}$$

$$x_{i+1} = x_i - \frac{(x_i^2 - \alpha)}{2x_i} \rightarrow x_{i+1} = \frac{1}{2} \left(x_i + \frac{\alpha}{x_i} \right)$$

Quando o método de Newton converge, geralmente o faz muito rapidamente, que caracteriza uma vantagem em relação à bissecção. No entanto, no método de Newton é necessário conhecer $f'(x)$ explicitamente. Em alguns casos isso não é possível. O próximo método contorna essa situação, a um maior custo computacional e consequentemente uma convergência mais lenta.

MÉTODO DA SECANTE

Utilizando uma aproximação para a derivada (conforme visto na seção 1.2), tal como,



$$f'(x) \approx \frac{f(x_i) - f(x_{i-1})}{x_i - x_{i-1}}$$

e substituindo na equação (4.6), obtemos:

$$x_{i+1} = x_i - f(x_i) \left[\frac{(x_i - x_{i-1})}{f(x_i) - f(x_{i-1})} \right], i = 1, 2, \dots \quad (4.7)$$

No cálculo de x_{i+1} necessitamos dos valores iniciais de x_i e x_{i-1} . No entanto, cada novo x_{i+1} requer apenas uma nova avaliação do valor de $f(x_i)$.

A interpretação gráfica do método da secante é similar ao de Newton. Ou seja, a reta tangente é substituída pela reta secante.

Exemplo 4.5

Um objeto é arremessado para cima com velocidade inicial $v_0 = 30m.s^{-1}$ a partir de uma altura $h_0 = 5m$, em um local onde a aceleração da gravidade é $g = -9,81m.s^{-2}$. Sabendo que

$$h(t) = h_0 + v_0t + \frac{gt^2}{2}$$

qual será o tempo gasto para o objeto tocar o solo, considerando o atrito desprezível?

$$h_0 + v_0t_i + \frac{g}{2}t_i^2 = 0 \rightarrow f(t_i) = 5 + 30t_i - 4,905t_i^2$$

Tabela 4.4 – Solução do exercício 4.6

i	t_i	t_{i-1}	$f(t_i)$	$f(t_{i-1})$	t_{i+1}	erro
1	8,0000	4,0000	-68,856	46,536	5,6131	—
2	5,6131	8,0000	18,8820	-68,856	6,1268	0,0838
3	6,1268	5,6131	4,7186	18,8820	6,2979	0,0271
4	6,2979	6,1268	-0,5747	4,7186	6,2793	0,0295
5	6,2793	6,2979	0,0139	-0,5747	6,2798	$6,9 \cdot 10^{-5}$



ITERAÇÃO FUNCIONAL E PONTO FIXO

O método de Newton é um exemplo onde uma sequência de pontos é calculada a partir de uma expressão da forma,

$$x_{i+1} = \phi(x_i), \quad (i \geq 0)$$

O algoritmo definido de tal forma é denominado de iteração funcional. No método de Newton, a função ϕ é dada por,

$$\phi(x_i) = x_i - \frac{f(x_i)}{f'(x_i)}.$$

A expressão (4.6) pode ser utilizada para gerar sequências que podem não convergir. No entanto, estamos interessados em casos onde $\lim_{i \rightarrow \infty} x_i$ existe. Supondo que,

$$\lim_{i \rightarrow \infty} x_i = \alpha.$$

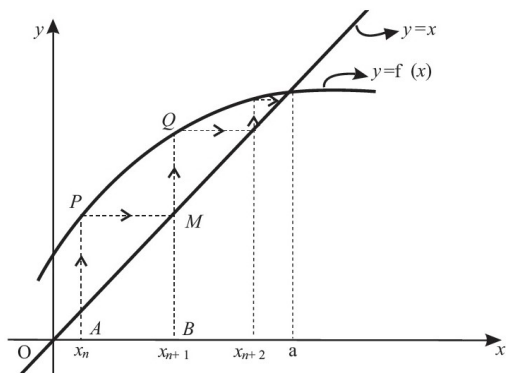
Se ϕ é contínua, então

$$\phi(\alpha) = \phi\left(\lim_{i \rightarrow \infty} x_i\right) = \lim_{i \rightarrow \infty} \phi(x_i) = \lim_{i \rightarrow \infty} x_{i+1} = \gamma \Rightarrow \phi(\alpha) = \gamma.$$

Onde γ é um **ponto fixo** da função ϕ . Intuitivamente, pensar num ponto fixo, é denominar um ponto onde a função “trava” no processo iterativo. Frequentemente, um problema matemático pode ser reduzido a um problema de encontrar um ponto fixo de uma função. Aplicações interessantes podem ser vistas em equações diferenciais, teoria da otimização e outras áreas.



Figura 4.6 – Exemplo do método iterativo funcional



Teorema 4.2

Seja $\varphi(x)$ contínua e diferenciável no intervalo $[a; b]$ tal que $\varphi([a; b]) \subset [a; b]$

$$\lambda = \max_{a \leq x \leq b} |\phi'(x)| < 1.$$

Então,

- $x = \varphi(x)$ possui uma única solução $\alpha \in [a; b]$;
- para alguma escolha de $x_0 \in [a; b]$, com $x_{i+1} = \Phi(x_i), i \geq 0$



$$\lim_{i \rightarrow \infty} x_i = \alpha$$

•

$$|\alpha - x_i| \leq \lambda^i |\alpha - x_0| \leq \frac{\lambda^i}{1 - \lambda} |x_1 - x_0|$$

$$\lim_{i \rightarrow \infty} \frac{\alpha - x_{i+1}}{\alpha - x_i} = \phi'(\alpha).$$

Ou seja, como condição suficiente, temos que se $|\phi'(x)| < 1$, para todo $x \in [a; b]$, então a sequência $(x_i)_{i \in \mathbb{N}}$ gerada pela expressão (4.8) converge para a raiz.

Exemplo 4.6

Utilizando o método de iteração funcional resolva a equação $\cos(x) - x = 0$ com $tol \leq 10^{-2}$.

Tabela 4.5 – Solução do exemplo 4.7

i	x_i	x_{i+1}	erro
0	0,5000	0,8775	—
1	0,8775	0,6390	0,3733
2	0,6390	0,8027	0,2039
3	0,8027	0,6947	0,1553
4	0,6947	0,7682	0,0955
5	0,7682	0,7191	0,0681
6	0,7191	0,7523	0,0441
7	0,7523	0,7301	0,0305
8	0,7301	0,7451	0,0201
9	0,7451	0,7350	0,0137
10	0,7350	0,7418	0,0091



Exemplo 4.7

Encontre opções da função $\phi(x)$ para o seguinte problema:

$$x^3 + x - 1 = 0$$

$$x^3 + x - 1 = 0 \rightarrow x^3 + x - x - 1 = 0 - x \rightarrow x = 1 - x^3 \rightarrow \phi(x) = 1 - x^3$$

ou

$$x^3 + x - 1 = 0 \rightarrow x^3 - x^3 + x - 1 = 0 - x^3 \rightarrow x = \sqrt[3]{1-x} \rightarrow \phi(x) = \sqrt[3]{1-x}$$

ou

$$x^3 + x - 1 = 0 \rightarrow x^3 + 2x^3 + x = 1 + 2x^3 \rightarrow x(3x^2 + 1) = 1 + 2x^2 \rightarrow x = \frac{1 + 2x^2}{1 + 3x^2} \rightarrow \phi(x) = \frac{1 + 2x^2}{1 + 3x^2}$$

SISTEMAS DE EQUAÇÕES NÃO LINEARES

Iremos considerar a problema de determinar as raízes do seguinte sistema de equações não lineares:

$$\begin{cases} f(x, y) = 0 \\ g(x, y) = 0 \end{cases} \quad (4.9)$$

Geometricamente, a solução desse sistema são os pontos no plano xy onde as curvas definidas por f e g se interceptam.

MÉTODO DE NEWTON NA SOLUÇÃO DE SISTEMAS NÃO LINEARES

Nesta subsecção iremos mostrar a construção de uma fórmula iterativa do método de Newton para uma função $F : \mathbb{R}^2 \rightarrow \mathbb{R}^2$. Para isso, Seja (x_0, y_0) uma aproximação para a solução do sistema definido na equação (4.9). Admitindo que f e g sejam suficientemente



diferenciáveis, expandimos $f(x, y)$ e $g(x, y)$, usando série de Taylor para funções de duas variáveis em torno de (x_0, y_0) . Assim:

$$\begin{cases} f(x, y) = f(x_0, y_0) + f_x(x_0, y_0)(x - x_0) + f_y(x_0, y_0)(y - y_0) + \dots \\ g(x, y) = g(x_0, y_0) + g_x(x_0, y_0)(x - x_0) + g_y(x_0, y_0)(y - y_0) + \dots \end{cases} \quad (4.10)$$

Admitindo que (x_0, y_0) esteja suficientemente próximo da solução a ponto de desprezarmos os termos de alta ordem da série em (4.10) e fazendo $f(x, y) = 0$ e $g(x, y) = 0$, teremos:

$$\begin{cases} f_x(x - x_0) + f_y(y - y_0) = -f \\ g_x(x - x_0) + g_y(y - y_0) = -g \end{cases} \implies \begin{cases} xf_x + yf_y = x_0f_x + y_0f_y - f \\ xg_x + yg_y = x_0g_x + y_0g_y - g \end{cases} \quad (4.11)$$

que pode ser escrito na seguinte forma matricial,

$$\begin{pmatrix} f_x & f_y \\ g_x & g_y \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} f_x & f_y \\ g_x & g_y \end{pmatrix} \begin{pmatrix} x_0 \\ y_0 \end{pmatrix} - \begin{pmatrix} f \\ g \end{pmatrix} \quad (4.12)$$

considerando que a matriz (Jacobiano) é invertível, temos a seguinte relação:

$$\begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} x_0 \\ y_0 \end{pmatrix} - \begin{pmatrix} f_x & f_y \\ g_x & g_y \end{pmatrix}^{-1} \begin{pmatrix} f \\ g \end{pmatrix} \quad (4.13)$$

todas as funções e derivadas parciais devem ser calculadas em (x_0, y_0) . Resolvendo a equação vetorial (4.13) de forma iterativa, obtemos:

$$\begin{cases} x_{i+1} = x_i - \left[\frac{f \cdot g_y - g \cdot f_y}{D(f, g)} \right]_{x_i, y_i} \\ y_{i+1} = y_i - \left[\frac{g \cdot f_x - f \cdot g_x}{D(f, g)} \right]_{x_i, y_i} \end{cases} \quad (4.14)$$

onde $D(f, g) = f_x \cdot g_y - f_y \cdot g_x$ representa o determinante do Jacobiano. De (4.11) também podemos tirar a seguinte relação:



$$\begin{pmatrix} f_x & f_y \\ g_x & g_y \end{pmatrix} \begin{pmatrix} \Delta x \\ \Delta y \end{pmatrix} = \begin{pmatrix} -f \\ -g \end{pmatrix}$$

onde $\Delta x = x - x_0$ e $\Delta y = y - y_0$. Que corresponde ao seguinte sistema de equações lineares,

$$\begin{cases} f_x \Delta x + f_y \Delta y = -f \\ g_x \Delta x + g_y \Delta y = -g \end{cases}$$

cuja solução será Δx e Δy . Assim,

$$\begin{cases} x_{i+1} = x_i + \Delta x_i \\ y_{i+1} = y_i + \Delta y_i \end{cases}$$

que corresponde a uma solução alternativa ao sistema de equações não lineares apresentado em (4.9).

Exemplo 4.8

Determinar uma raiz do sistema não linear:

$$\begin{cases} x^2 + y^2 = 2 \\ x^2 - y^2 = 1 \end{cases} \quad (4.15)$$

com tolerância de 10^{-3} , usando o método de Newton e considerando $(x_0; y_0) = (1, 2; 0, 7)$

$$\begin{cases} f(x, y) = x^2 + y^2 - 2 \\ g(x, y) = x^2 - y^2 - 1 \end{cases}$$

com

$$\begin{pmatrix} f_x & f_y \\ g_x & g_y \end{pmatrix} = \begin{pmatrix} 2x & 2y \\ 2x & -2y \end{pmatrix} \xrightarrow{x_0 = 1, 2; y_0 = 0, 7} \begin{pmatrix} 2,4 & 1,4 \\ 2,4 & -1,4 \end{pmatrix}$$

e, $D = -6,72$, $f(1, 2; 0, 7) = -0,07$, $g(1, 2; 0, 7) = -0,05$. Assim,



$$X_1 = \begin{pmatrix} x_1 \\ y_1 \end{pmatrix} = \begin{pmatrix} 1,2 \\ 0,7 \end{pmatrix} + \frac{1}{6,72} \begin{pmatrix} -1,4 & -1,4 \\ -2,4 & 2,4 \end{pmatrix} \begin{pmatrix} 0,07 \\ 0,05 \end{pmatrix} = \begin{pmatrix} 1,175 \\ 0,6928 \end{pmatrix}$$

$$X_2 = \begin{pmatrix} x_2 \\ y_2 \end{pmatrix} = \begin{pmatrix} 1,175 \\ 0,6928 \end{pmatrix} + \frac{1}{6,5123} \begin{pmatrix} -1,3856 & -1,3856 \\ -2,35 & 2,35 \end{pmatrix} \begin{pmatrix} 0,1394 \\ 0,05 \end{pmatrix} = \begin{pmatrix} 1,1242 \\ 1,1605 \end{pmatrix}$$

com erro relativo de,

$$E_r = \frac{\|X_2 - X_1\|_\infty}{\|X_2\|_\infty} \approx 0,4030$$

Algorithm 3 Método de Newton para Sistema de Equações Não lineares

Passo 1: $k = 0$.

Passo 2: Calcule $\mathbf{F}(\mathbf{x}_k)$.

if $\|\mathbf{F}(\mathbf{x}_k)\| \leq tol$ **then**

$\mathbf{x}^* = \mathbf{x}_k$

else

 Passo 3: Calcule $\mathbf{J}(\mathbf{x}_k)$

end if

Passo 4: obtenha Δ , solução do sistema linear: $\mathbf{J}(\mathbf{x}_k)\Delta = -\mathbf{F}(\mathbf{x}_k)$

Passo 5: faça $\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta_k$

if $\|\mathbf{x}_{k+1} - \mathbf{x}_k\| \leq tol$ **then**

\mathbf{x}_{k+1} é solução

else

$k = k + 1$ e volte ao passo 2.

end if



CONVERGÊNCIA

Teorema 4.3

Suponha que $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ é continuamente diferenciável e $F(\mathbf{x}^*) = \mathbf{0}$. Se,

- A matriz Jacobiana $J(\mathbf{x}^*)$ de F em \mathbf{x}^* é não singular ($|J(\mathbf{x}^*)| \neq 0$).
- J é Lipschitz numa vizinhança de \mathbf{x}^* , então para todo $b f x_0$ suficientemente próxima de \mathbf{x}^* , o método de Newton produz uma sequência $\{\mathbf{x}_k\}$ que converge quadraticamente para \mathbf{x}^* .

Exemplo 4.9

Resolver o sistema de equações não lineares, de modo que a matriz $F(\mathbf{x})$ é dada por,

$$\begin{pmatrix} x^2 - y - 1 \\ xy - 1 \end{pmatrix}$$

com condição inicial $\mathbf{x}_0 = (2, 1)^T$ e $tol < 10^{-4}$.

Tabela 4.6 – Solução do exemplo 4.10.

k	\mathbf{x}_k	$\ \mathbf{x}_k - \mathbf{x}_{k-1}\ $
0	$(2, 1)^T$	
1	$(1,444444, 0,777777)^T$	0,55
2	$(1,3294, 0,7542)^T$	0,11
3	$(1,324724, 0,754871)^T$	0,0047
4	$(1,324717, 0,754877)^T$	0,0000062



Exemplo 4.10

Aplicação.

Cada vez que um GPS é usado, o sistema de equações não lineares

$$\begin{cases} (x - a_1)^2 + (y - b_1)^2 + (z - c_1)^2 = [(C(t_1 - D))]^2 \\ (x - a_2)^2 + (y - b_2)^2 + (z - c_2)^2 = [(C(t_2 - D))]^2 \\ (x - a_3)^2 + (y - b_3)^2 + (z - c_3)^2 = [(C(t_3 - D))]^2 \\ (x - a_4)^2 + (y - b_4)^2 + (z - c_4)^2 = [(C(t_4 - D))]^2 \end{cases}$$

é resolvido para as coordenadas (x, y, z) do receptor. Para cada satélite i as localizações são (a_i, b_i, c_i) , e t_i é o tempo de sincronização da transmissão a partir do satélite. Além disso, temos a constante de velocidade da luz C , D que representa a diferença entre o tempo de sincronização do satélite e do receptor.

EXERCÍCIOS PROPOSTOS

Exercício 4.5

Resolva a seguinte equação,

$$e^x = x^e$$

Exercício 4.6

Utilizando o método de Newton, mostrar que

$$\sqrt[p]{a} \approx x_{i+1} = \frac{1}{p} \left[(p-1)x_i + \frac{a}{x_i^{p-1}} \right]$$



Exercício 4.7

Determinar

$$\sqrt[8]{3}$$

utilizando o método de sua preferência com $E_r \leq 10^{-3}$.

Exercício 4.8

A localização \bar{x} do centróide de um setor circular é dada por:

$$\bar{x} = \frac{2r \operatorname{sen}(\theta)}{3\theta}.$$

Determine o ângulo θ para o qual $\bar{x} = r$. Utilizando os seguintes métodos:

- Método da bissecção com $a = 1$ e $b = 2$. Realize as cinco primeiras iterações;
- Método da secante e comece com $x_0 = 1$ e $x_1 = 2$. Calcule as cinco primeiras iterações;
- Método de Newton. Comece em $x_0 = 1$ e realize as cinco primeiras iterações.



Exercício 4.9

A área S da superfície lateral de um cone é dada por:

$$S = \pi r \sqrt{r^2 + h^2}$$

onde r é o raio da base e h a altura. Determine o raio de um cone que tenha uma área superficial de 1200 m^2 e uma altura de 20 m , calculando cinco iterações com o método da iteração de ponto fixo. Comece com $r_0 = 17 \text{ m}$.

Exercício 4.10

Usando o método de Newton determine $x \in \mathbb{R}$, com $E_a \leq 10^{-2}$, tal que a matriz:

$$\begin{pmatrix} 0.5 & 0.2 & x \\ 0.4 & x & 0.5 \\ x & 0.5 & 0.2 \end{pmatrix}$$

seja singular.

Exercício 4.11

Encontre uma raiz da equação,



$$f(x) = x^3 - 9x^2 + 25x \left(1 + \frac{\sin^2(x)}{25} \right) + \frac{x}{\sec^2(x)} - 24$$

no intervalo $1 \leq x \leq 2,5$ com $Er \leq 10^{-3}$.

Exercício 4.12

Num estudo de coleta em energia solar localizado num campo de espelhos planos, com um coletor central, um pesquisador obteve a seguinte expressão para o fator de concentração geométrico G :

$$G = \frac{\pi(h/\cos(\theta))^2 F}{0,5\pi D^2(1 + \sin(\theta) - 0,5\cos(\theta))}$$

onde θ é o ângulo de borda do campo, F é a cobertura fracionária do campo de espelhos, D é o diâmetro do coletor e, h é a altura do coletor. Determine o ângulo θ se, $h = 300$, $G = 1200$, $F = 0,8$ e $D = 14$.

Exercício 4.13

A equação de Kepler para determinar órbitas de satélites, é dada por

$$M = x - E \cdot \sin(x).$$



Dado que $E = 0, 2$ e $M = 0, 5$, obtenha a raiz da equação de Kepler usando o método de Newton.

Exercício 4.14

A velocidade ascendente de um foguete pode ser calculada pela seguinte equação:

$$v = u \cdot \ln \left(\frac{m_0}{m_0 - q \cdot t} \right) - g \cdot t$$

onde v representa a velocidade de subida, u é a velocidade na qual o combustível é repelido com relação ao foguete, m_0 é a massa inicial ($t = 0s$), q é a taxa de consumo de combustível, e g a aceleração da gravidade para baixo $g \approx 9,81m/s^2$. Se $u = 2000m/s$, $m_0 = 150.000kg$, e $q = 2700kg/s$. Calcule o instante no qual a velocidade é igual a $750m/s$ com $10 \leq t \leq 50s$. Considere uma tolerância de 10^{-2} .

Exercício 4.15

Um tanque de armazenamento contém um líquido à profundidade y . É tirado líquido a uma vazão constante Q , para atender a demanda, o conteúdo é repostado a uma taxa senoidal de $3Qsen^2(t)$. A equação que descreve o sistema é dada por,

$$A \cdot dy = (3Qsen^2(t) - Q)dt.$$



Use um dos métodos estudados para encontrar o instante de tempo em que a profundidade $y = 1$ m, considerando $t \in [0 ; 10]$. Suponha $A = 1200$ m , $Q = 500$ e $E_a \leq 10^{-4}$.



UNIDADE II



SISTEMAS DE EQUAÇÕES LINEARES

Em muitas aplicações práticas nas ciências e nas engenharias, os dados são frequentemente organizadas em linhas e colunas formando um sistema de equações lineares e para extrair informações relevantes, de considerável importância, que são obtidas por meio da resolução de tais sistemas. O objetivo principal deste capítulo é o estudo de técnicas de resolução por meio de algoritmos computacionais dos sistemas de equações lineares - SEL.

Uma equação linear nas incógnitas x_1, x_2, \dots, x_n é uma equação que pode ser escrita como:

$$a_1 x_1 + a_2 x_2 + \dots + a_n x_n = b$$

na qual as variáveis a_1, a_2, \dots, a_n são constantes, e b é o termo independente. Os valores das incógnitas que transformam uma equação linear em identidade, ou seja, satisfazem a equação, constituem uma solução.

Exemplo 5.1

Seja a equação $2x_1 + 3x_2 = 18$. Podemos observar que $x_1 = 3$ e $x_2 = 4$. É solução da equação. Ou seja,

$$2x_1 + 3x_2 = 2 \cdot 3 + 3 \cdot 4 = 18.$$



No entanto, se pensarmos em espaços vetoriais, podemos definir uma classe de mapeamentos $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$, onde o i -ésimo componente de $\mathbf{y} = T(\mathbf{x})$ é expresso em termos dos componentes x_j do vetor \mathbf{x} por meio de,

$$y_i = \sum_{j=1}^n a_{ij}x_j, \quad i = 1, 2, \dots, m \quad (5.1)$$

onde os a_{ij} são escalares. Esses mapeamentos são lineares, consequentemente, cada mapeamento linear

$$T : \mathbb{R}^n \rightarrow \mathbb{R}^m \\ T(\mathbf{x}) \mapsto \mathbf{y}$$

pode ser escrito da forma 5.1. Ou seja, o vetor $\mathbf{x} \in \mathbb{R}^n$ pode ser expresso como uma combinação linear de vetores da base canônica, $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n$, onde \mathbf{e}_j possui o j -ésimo componente unitário, e os demais nulos.

$$\mathbf{x} = \sum_{j=1}^n x_j \mathbf{e}_j.$$

Sendo T linear,

$$\mathbf{y} = T(\mathbf{x}) = \sum_{j=1}^n x_j T(\mathbf{e}_j).$$

A um conjunto com m **equações lineares** e n **incógnitas**, denominamos de Sistemas de Equações Lineares. Como exemplo,



$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2 \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n = b_m. \end{cases}$$

onde $a_{ij}, b_j, 1 \leq i \leq m, 1 \leq j \leq n$, são números reais (ou complexos) conhecidos. Se $b_i = 0, \forall i$ dizemos que o sistema é homogêneo. Caso contrário, denominamos de não homogêneo. Uma solução do SEL é uma n -upla (x_1, x_2, \dots, x_n) que satisfaz simultaneamente as m **equações lineares**. O sistema pode ser escrito numa forma matricial. Ou seja,

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \ddots & \cdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}_{m \times n} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}_{n \times 1} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix}_{m \times 1}$$

$$\mathbf{Ax} = \mathbf{b}$$

sendo,

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \ddots & \cdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix} \quad \text{e} \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$



de modo que \mathbf{A} representa a matriz dos coeficientes de dimensão $m \times n$, \mathbf{b} o vetor coluna de parâmetros conhecidos ($m \times 1$) e \mathbf{x} o vetor coluna de parâmetros desconhecidos ($m \times 1$).

Uma outra matriz que pode ser associada ao sistema é da forma,

$$\left(\begin{array}{cccc|c} a_{11} & a_{12} & \cdots & a_{1n} & b_1 \\ a_{21} & a_{22} & \cdots & a_{2n} & b_2 \\ \vdots & \ddots & \cdots & \vdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} & b_m \end{array} \right)$$

$$[\mathbf{A}|\mathbf{b}].$$

denominada simplesmente de **matriz ampliada ou aumentada** do sistema, onde cada linha nesta matriz corresponde a uma representação simplificada da equação correspondente ao SEL original. Se $m = n$ temos simplesmente,

$$\left(\begin{array}{cccc} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \ddots & \cdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{nn} \end{array} \right) \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}$$

$$\mathbf{Ax} = \mathbf{b}$$

Neste caso, se a matriz quadrada \mathbf{A} for invertível, a resolução do SEL é dada por meio da seguinte equação matricial:

$$\mathbf{x} = \mathbf{A}^{-1} \cdot \mathbf{b}.$$

No entanto, em muitas aplicações práticas, resolver um SEL baseando-se no conceito da matriz inversa, mesmo sendo um método eficiente, não é prático e possui um custo computacional alto.



Portanto, um procedimento prático são os métodos em que obtemos um sistema equivalente ao sistema dado por meio da aplicação de uma sequência de operações elementares. De modo que, sistemas lineares equivalentes possuem o mesmo conjunto solução. As principais operações elementares são:

- Permutar duas equações;
- Multiplicar uma equação por uma constante não nula;
- Somar ou subtrair equações entre si.

Exemplo 5.2

Seja o SEL,

$$\begin{cases} 2x_1 + 3x_2 = 18 \\ 3x_1 + 4x_2 = 25 \end{cases}$$

Que pode ser escrito na seguinte forma matricial:

$$\begin{pmatrix} 2 & 3 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 18 \\ 25 \end{pmatrix} \iff \begin{pmatrix} 2 & 3 & | & 18 \\ 3 & 4 & | & 25 \end{pmatrix}$$

Assim podemos efetuar as operações elementares na matriz equivalente ao SEL original. Para isso, iremos começar a estudar alguns métodos, começando pelos métodos diretos, em que se destacam o método de eliminação de Gauss e o de Gauss- Jordan, que elimina partes da matriz aumentada por meio de operações elementares nos vetores linha, de modo a obter um sistema triangular equivalente. bem como os métodos de decomposição, que são um caso particular do métodos diretos. Em seguida, estudaremos os chamados métodos iterativos.



MÉTODOS DIRETOS

MÉTODO DE ELIMINAÇÃO DE GAUSS

A resolução de SEL utilizando o método de Eliminação de Gauss, consiste primeiramente em escrever o sistema original na forma de uma **matriz aumentada**, que consiste na matriz de coeficientes **A** concatenada com o vetor solução **b**. Ou seja,

$$\left(\begin{array}{cccc|c} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} & b_1 \\ a_{21} & a_{22} & a_{23} & \cdots & a_{2n} & b_2 \\ \vdots & \ddots & \cdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \cdots & a_{nn} & b_n \end{array} \right)$$

$$[\mathbf{A}|\mathbf{b}].$$

O método de eliminação de Gauss consiste em efetuar operações elementares nas linhas da matriz aumentada, tomando os elementos da diagonal como pivô, de modo a transformar a matriz no equivalente a um SEL triangular superior (ou inferior), tomando em cada passo o elemento da diagonal principal como pivô.

$$[\mathbf{A}|\mathbf{b}] \xrightarrow{(op.\ element.)} [\mathbf{A}^*|\mathbf{b}^*]$$

Exemplo 5.3

Seja o SEL:

$$\left(\begin{array}{ccc|c} a_{11} & a_{12} & a_{13} & b_1 \\ a_{21} & a_{22} & a_{23} & b_2 \\ a_{31} & a_{32} & a_{33} & b_3 \end{array} \right)$$



que após a eliminação progressiva, o SEL se resume a,

$$\left(\begin{array}{ccc|c} a_{11} & a_{12} & a_{13} & b_1 \\ 0 & a'_{22} & a'_{23} & b'_2 \\ 0 & 0 & a'_{33} & b'_3 \end{array} \right)$$

podendo ser resolvido via aplicação da substituição regressiva apresentando, em seguida, o vetor solução.

Exemplo 5.4

Use o método de eliminação de Gauss para resolver o seguinte SEL:

$$\begin{cases} 2x_1 + 3x_2 - x_3 = 5 \\ 4x_1 + 4x_2 - 3x_3 = 3 \\ 2x_1 - 3x_2 + x_3 = -1. \end{cases}$$

Etapa 01 - Escrever a matriz aumentada do sistema e escolha do pivô ($\neq 0$), ou seja

$$\left(\begin{array}{ccc|c} 2 & 3 & -1 & 5 \\ 4 & 4 & -3 & 3 \\ 2 & -3 & 1 & -1 \end{array} \right)$$

Etapa 02 - Operações elementares sobre as linhas da matriz aumentada até obter uma matriz triangular superior (ou inferior).

$$\left(\begin{array}{ccc|c} 2 & 3 & -1 & 5 \\ 0 & -2 & -1 & -7 \\ 0 & 0 & 5 & 15 \end{array} \right)$$



Etapa 03 - Aplicando as substituições retroativas, obtemos:

$$x_3 = \frac{15}{3} = 3$$

$$x_2 = \frac{-7+3}{2} = 2$$

$$x_1 = \frac{5-3 \cdot 2+3}{2} = 1$$

Armadilhas dos métodos de eliminação de Gauss

- Divisão por zero : tanto durante a fase de eliminação quanto durante a de substituição é possível ocorrer uma divisão por zero.

$$\begin{cases} 0x_1 + 2x_2 + 3x_3 = 8 \\ 4x_1 + 6x_2 + 7x_3 = -3 \\ 2x_1 + x_2 + 6x_3 = 5 \end{cases}$$

Também podem surgir problemas quando um coeficiente está muito próximo de zero. A técnica de pivotamento foi desenvolvida para evitar parcialmente esses problemas.

- Erros de arredondamento : O problema de erros de arredondamento pode tornar-se particularmente importante quando se resolve um número grande de equações, por causa do fato de que cada resultado depende dos resultados anteriores. Conseqüente- mente, um erro nas etapas iniciais tende a se propagar — isto é, irá causar erros nas etapas subseqüentes
- Sistemas mal condicionados : são aqueles para os quais pequenas mudanças nos coeficientes resultam em grandes mudanças nas soluções.

Técnicas para melhorar a solução

- Uso de mais variáveis significativas



- Pivotamento parcial: determinar o maior coeficiente disponível na coluna abaixo do elemento pivô. As linhas podem ser trocadas de modo que o maior coeficiente seja o elemento pivô.

MÉTODO DE GAUSS - JORDAN

O método de Gauss - Jordan, consiste numa modificação no método de Gauss, de modo a transformar o SEL, num sistema equivalente com a matriz dos coeficientes sendo a matriz identidade.

$$[\mathbf{A}|\mathbf{b}] \xrightarrow{\text{(op. element.)}} [\mathbf{I}|\mathbf{b}^*]$$

uma outra aplicação do método de Gauss - Jordan é o de encontrar a matriz inversa, consistindo basicamente na resolução do SEL,

$$\mathbf{AX} = \mathbf{I} \rightarrow \mathbf{X} = \mathbf{A}^{-1}$$

Exemplo 5.5

Resolva o SEL, utilizando o método de eliminação de Gauss - Jordan.

$$\begin{cases} 3x_1 + 2x_2 + 4x_3 = 1 \\ x_1 + x_2 + 2x_3 = 2 \\ 4x_1 + 3x_2 - 2x_3 = 3 \end{cases}$$

e encontre a matriz inversa de,

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 2 \\ 0 & 1 & 1 \end{pmatrix}$$



MÉTODOS DE DECOMPOSIÇÃO

Consiste em decompor a matriz \mathbf{A} no produto de uma matriz triangular inferior \mathbf{L} e uma matriz triangular superior \mathbf{U} . Com isso podemos utilizar o método de decomposição na resolução de sistemas $\mathbf{Ax} = \mathbf{b}$, bem como o cálculo da inversa da matriz \mathbf{A} . É um método vantajoso quando se necessita resolver um sistema de equações para inúmeros vetores \mathbf{b} 's diferentes.

Outros métodos de decomposição consistem em,

- Decomposição *Doolittle*: \mathbf{L} possui uma diagonal unitária;
- Decomposição *CROUT*: \mathbf{U} possui uma diagonal unitária;
- Decomposição *Cholesky*: $\mathbf{U} = \mathbf{L}^T$ ou $\mathbf{L} = \mathbf{U}^T$. Em particular,

$$l_{ii} = u_{ii} \text{ para } i = 1, 2, \dots, n.$$

Exemplo 5.6

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{pmatrix} \cdot \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{pmatrix}$$

$$\begin{cases} 3x_1 - 0,1x_2 - 0,2x_3 = 7,85 \\ 0,1x_1 + 7x_2 - 0,3x_3 = -19,3 \\ 0,3x_1 - 0,2x_2 + 10x_3 = 71,4 \end{cases}$$

DECOMPOSIÇÃO CHOLESKY

Uma matriz simétrica $\mathbf{A} = \mathbf{A}^T$ e definida positiva $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$ pode ser decomposta no produto de uma matriz triangular inferior e sua matriz adjunta, o que é bastante útil, por exemplo, para a solução numérica eficiente e simulações de Monte Carlo. Ou seja,



$$\mathbf{A} = \mathbf{G}\mathbf{G}^T$$

onde a matriz triangular inferior \mathbf{G} é denominada de fator *Cholesky* da matriz \mathbf{A} . A decomposição de *Cholesky* é usada principalmente na resolução de sistemas de equações lineares, por meio das seguintes relações,

$$\mathbf{A} = \mathbf{G}\mathbf{G}^T \Rightarrow \begin{cases} \mathbf{G}\mathbf{d} = \mathbf{b} \\ \mathbf{G}^T \mathbf{x} = \mathbf{d} \end{cases}$$

Para isso, iniciamos decompondo a matriz do sistema. Por exemplo,

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = \begin{pmatrix} g_{11} & 0 & 0 \\ g_{21} & g_{22} & 0 \\ g_{31} & g_{32} & g_{33} \end{pmatrix} \cdot \begin{pmatrix} g_{11} & g_{12} & g_{13} \\ 0 & g_{22} & g_{23} \\ 0 & 0 & g_{33} \end{pmatrix}$$

ou seja,

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = \begin{pmatrix} g_{11}^2 & g_{11}\cdot g_{12} & g_{11}\cdot g_{13} \\ g_{21}\cdot g_{11} & g_{21}^2\cdot g_{22}^2 & g_{21}\cdot g_{13} + g_{22}\cdot g_{23} \\ g_{31}\cdot g_{11} & g_{31}\cdot g_{12} + g_{32}\cdot g_{22} & g_{31}^2 + g_{32}^2 + g_{33}^2 \end{pmatrix}$$

assim, podemos tirar as seguintes relações:

$$\left\{ \begin{array}{l} a_{11} = g_{11}^2 \rightarrow g_{11} = \sqrt{a_{11}} \\ a_{21} = g_{21}\cdot g_{11} \rightarrow g_{21} = a_{21}/g_{11} \\ a_{22} = g_{21}^2\cdot g_{22}^2 \rightarrow g_{22} = \sqrt{a_{22} - g_{21}^2} \\ a_{31} = g_{31}\cdot g_{11} \rightarrow g_{31} = a_{31}/g_{11} \\ a_{32} = g_{31}\cdot g_{12} + g_{32}\cdot g_{22} \rightarrow g_{32} = (a_{32} - g_{31}\cdot g_{12})/g_{22} \\ a_{33} = g_{31}^2 + g_{32}^2 + g_{33}^2 \rightarrow g_{33} = \sqrt{a_{33} - (g_{31}^2 + g_{32}^2)} \end{array} \right.$$



generalizando, temos:

$$g_{jj} = (\pm) \sqrt{a_{jj} - \sum_{k=1}^{j-1} g_{jk}^2}$$

$$g_{ij} = \frac{1}{g_{jj}} \left(a_{ij} - \sum_{k=1}^{j-1} g_{ik} \cdot g_{jk} \right), \text{ para } i > j.$$

Exemplo 5.7

Para a seguinte matriz **A**, encontrar a matriz de decomposição de Cholesky **G**.

$$A = \begin{pmatrix} 25 & 15 & -5 \\ 15 & 18 & 0 \\ -5 & 0 & 11 \end{pmatrix}$$

$$\left\{ \begin{array}{l} g_{11} = \sqrt{25} = 5 \\ g_{21} = 15/5 = 3 \\ g_{22} = \sqrt{18 - 9} = 3 \\ g_{31} = -5/5 = -1 \\ g_{32} = (0 - 3 \cdot (-1))/3 = 1 \\ g_{33} = \sqrt{11 - ((-1)^2 + 1^2)} = 3 \end{array} \right.$$

$$G = \begin{pmatrix} 5 & 0 & 0 \\ 3 & 3 & 0 \\ -1 & 1 & 3 \end{pmatrix}$$



DECOMPOSIÇÃO LU

Dado um sistema $\mathbf{A}\mathbf{x} = \mathbf{b}$ com $\mathbf{A} \in \mathbb{R}^{n \times n}$, cujos menores são não singulares. Assim, o método de decomposição LU consiste em fatorar a matriz \mathbf{A} como um produto de duas matrizes, $\mathbf{A} = \mathbf{L}\mathbf{U}$, de modo que uma seja **triangular inferior**, \mathbf{L} e a outra **triangular superior**, \mathbf{U} . O sistema resultante a ser resolvido se resume a,

$$\mathbf{Ax} = \mathbf{b} \rightarrow \mathbf{L}\mathbf{U} = \mathbf{b} \rightarrow \underbrace{\mathbf{L}\mathbf{U}}_{\mathbf{y}}\mathbf{x} = \mathbf{b}.$$

Ou seja, O nosso objetivo é decompor a matriz \mathbf{A} em duas matrizes, uma triangular inferior (\mathbf{L}) e uma triangular superior (\mathbf{U}).

$$\underbrace{\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}}_{\mathbf{A}} = \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{pmatrix}}_{\mathbf{L}} \cdot \underbrace{\begin{pmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{pmatrix}}_{\mathbf{U}}$$

$$\mathbf{A} = \mathbf{L}\mathbf{U}$$

E, em seguida resolver por meio da substituição progressiva,

$$\mathbf{Ly} = \mathbf{b}$$

e por substituição regressiva,

$$\mathbf{Ux} = \mathbf{y}.$$



Exemplo da decomposição LU

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ l_{21}u_{11} & l_{21}u_{12} + u_{22} & l_{21}u_{13} + u_{23} \\ l_{31}u_{11} & l_{31}u_{12} + l_{32}u_{22} & l_{31}u_{13} + l_{32}u_{23} + u_{33} \end{bmatrix}$$

Daí, podemos concluir que:

$$a_{1n} = u_{1n}, n = 1, 2, 3$$

$$a_{21} = l_{21}u_{11} \rightarrow l_{21} = \frac{a_{21}}{u_{11}}$$

$$a_{22} = l_{21}u_{12} + u_{22} \rightarrow u_{22} = a_{22} - l_{21}u_{12}$$

$$a_{23} = l_{21}u_{13} + u_{23} \rightarrow u_{23} = a_{23} - l_{21}u_{13}$$

e,

$$a_{31} = l_{31}u_{11} \rightarrow l_{31} = \frac{a_{31}}{u_{11}}$$

$$a_{32} = l_{31}u_{12} + l_{32}u_{22} \rightarrow l_{32} = \frac{(a_{32} - l_{31}u_{12})}{u_{22}}$$

$$a_{33} = l_{31}u_{13} + l_{32}u_{23} + u_{33} \rightarrow u_{33} = (a_{33} - l_{31}u_{13}) - l_{32}u_{23}$$

DECOMPOSIÇÃO LU VIA MÉTODO DE ELIMINAÇÃO DE GAUSS

Podemos unir de maneira direta os métodos de eliminação de Gauss e decomposição LU, de modo a obter de maneira direta e



conjunta a solução do SEL, para isso iremos mostrar o processo por meio do seguinte exemplo. Considerando,

$$\mathbf{A}_a = \left(\begin{array}{cccc|c} a_{11} & a_{12} & a_{13} & a_{14} & b_1 \\ a_{21} & a_{22} & a_{23} & a_{24} & b_2 \\ a_{31} & a_{32} & a_{33} & a_{34} & b_3 \\ a_{41} & a_{42} & a_{43} & a_{44} & b_4 \end{array} \right)$$

Passo 1 - Escolha o pivô e dos multiplicadores : $p_1 \neq 0$. Geralmente $p_1 = a_{11}$ e os multiplicadores:

$$m_1 = \frac{a_{21}}{p_1}, m_2 = \frac{a_{31}}{p_1} \text{ e } m_3 = \frac{a_{41}}{p_1},$$

que correspondem, respectivamente aos elementos l_{21} , l_{31} e l_{41} da matriz \mathbf{L} .

Passo 2 - Calcular os termos da seguinte matriz, a partir da matriz original aumentada. Ou seja,

$$a_{32}^{(1)} = \frac{1}{p_1} \left| \begin{array}{cc|c} a_{11} & a_{12} & \\ a_{31} & a_{32} & \end{array} \right|, a_{33}^{(1)} = \frac{1}{p_1} \left| \begin{array}{cc|c} a_{11} & a_{13} & \\ a_{31} & a_{33} & \end{array} \right|, a_{34}^{(1)} = \frac{1}{p_1} \left| \begin{array}{cc|c} a_{11} & a_{14} & \\ a_{31} & a_{34} & \end{array} \right|, b_3^{(1)} = \frac{1}{p_1} \left| \begin{array}{cc|c} a_{11} & b_2 & \\ a_{31} & b_3 & \end{array} \right|$$

$$a_{42}^{(1)} = \frac{1}{p_1} \left| \begin{array}{cc|c} a_{11} & a_{12} & \\ a_{41} & a_{42} & \end{array} \right|, a_{43}^{(1)} = \frac{1}{p_1} \left| \begin{array}{cc|c} a_{11} & a_{13} & \\ a_{41} & a_{43} & \end{array} \right|, a_{44}^{(1)} = \frac{1}{p_1} \left| \begin{array}{cc|c} a_{11} & a_{14} & \\ a_{41} & a_{44} & \end{array} \right|, b_4^{(1)} = \frac{1}{p_1} \left| \begin{array}{cc|c} a_{11} & b_3 & \\ a_{31} & b_4 & \end{array} \right|$$

gerando a seguinte matriz,

$$\mathbf{A}_a^{(1)} = \left(\begin{array}{cccc|c} a_{11} & a_{12} & a_{13} & a_{14} & b_1 \\ 0 & a_{22}^{(1)} & a_{23}^{(1)} & a_{24}^{(1)} & b_2^{(1)} \\ 0 & a_{32}^{(1)} & a_{33}^{(1)} & a_{34}^{(1)} & b_3^{(1)} \\ 0 & a_{42}^{(1)} & a_{43}^{(1)} & a_{44}^{(1)} & b_4^{(1)} \end{array} \right)$$



escolhendo-se um novo pivô

$$p_2 = a_{22}^{(1)} \neq 0$$

e calculando os multiplicadores

$$m_4 = \frac{a_{32}^{(1)}}{p_2}$$

e

$$m_5 = \frac{a_{42}^{(1)}}{p_2},$$

correspondendo aos elementos l_{32} e l_{42} da matriz \mathbf{L} e, calculando-se os seguintes termos,

$$a_{33}^{(2)} = \frac{1}{p_2} \begin{vmatrix} a_{22}^{(1)} & a_{23}^{(1)} \\ a_{31}^{(1)} & a_{33}^{(1)} \end{vmatrix}, \quad a_{34}^{(2)} = \frac{1}{p_2} \begin{vmatrix} a_{22}^{(1)} & a_{24}^{(1)} \\ a_{31}^{(1)} & a_{34}^{(1)} \end{vmatrix}, \quad b_3^{(1)} = \frac{1}{p_2} \begin{vmatrix} a_{22}^{(1)} & b_2^{(1)} \\ a_{31}^{(1)} & b_3^{(1)} \end{vmatrix}$$

$$a_{43}^{(2)} = \frac{1}{p_2} \begin{vmatrix} a_{22}^{(1)} & a_{23}^{(1)} \\ a_{42}^{(1)} & a_{43}^{(1)} \end{vmatrix}, \quad a_{44}^{(2)} = \frac{1}{p_2} \begin{vmatrix} a_{22}^{(1)} & a_{24}^{(1)} \\ a_{42}^{(1)} & a_{44}^{(1)} \end{vmatrix}, \quad b_4^{(1)} = \frac{1}{p_2} \begin{vmatrix} a_{22}^{(1)} & b_3^{(1)} \\ a_{31}^{(1)} & b_4^{(1)} \end{vmatrix}$$

obteremos a seguinte matriz,

$$\mathbf{A}_a^{(2)} = \begin{pmatrix} a_{11} & a_{12} & a_{13} & a_{14} & b_1 \\ 0 & a_{22}^{(1)} & a_{23}^{(1)} & a_{24}^{(1)} & b_2^{(1)} \\ 0 & 0 & a_{33}^{(2)} & a_{34}^{(2)} & b_3^{(2)} \\ 0 & 0 & a_{43}^{(2)} & a_{44}^{(2)} & b_4^{(2)} \end{pmatrix}$$

e, finalmente escolhendo

$$p_3 = a_{33}^{(2)} \neq 0 \text{ e } m_6 = \frac{a_{43}^{(2)}}{p_3} = l_{43},$$



temos:

$$a_{44}^{(3)} = \frac{1}{p_3} \left| \begin{array}{cc} a_{33}^{(2)} & a_{34}^{(2)} \\ a_{43}^{(2)} & a_{44}^{(2)} \end{array} \right|, \quad b_4^{(3)} = \frac{1}{p_3} \left| \begin{array}{cc} a_{33}^{(1)} & b_3^{(2)} \\ a_{43}^{(1)} & b_4^{(2)} \end{array} \right|$$

ou seja,

$$\mathbf{A}_a^{(3)} = \left(\begin{array}{cccc|c} a_{11} & a_{12} & a_{13} & a_{14} & b_1 \\ 0 & a_{22}^{(1)} & a_{23}^{(1)} & a_{24}^{(1)} & b_2^{(1)} \\ 0 & 0 & a_{33}^{(2)} & a_{34}^{(2)} & b_3^{(2)} \\ 0 & 0 & 0 & a_{44}^{(3)} & b_4^{(3)} \end{array} \right)$$

que corresponde a matriz da eliminação de Gauss e de maneira similar, temos que

$$\mathbf{U} = \left(\begin{array}{cccc} a_{11} & a_{12} & a_{13} & a_{14} \\ 0 & a_{22}^{(1)} & a_{23}^{(1)} & a_{24}^{(1)} \\ 0 & 0 & a_{33}^{(2)} & a_{34}^{(2)} \\ 0 & 0 & 0 & a_{44}^{(3)} \end{array} \right), \quad \mathbf{L} = \left(\begin{array}{cccc} 1 & 0 & 0 & 0 \\ m_1 & 1 & 0 & 0 \\ m_2 & m_4 & 1 & 0 \\ m_3 & m_5 & m_6 & 1 \end{array} \right)$$

correspondendo as matrizes da decomposição LU.

MATRIZ INVERSA COM DECOMPOSIÇÃO LU

O método da eliminação de Gauss pode ser aplicado de maneira direta à matriz aumentada

$$[\mathbf{A}|\mathbf{I}]$$

como $\mathbf{A} = \mathbf{LU}$ e por meio da relação $\mathbf{AA}^{-1} = \mathbf{I} \rightarrow \mathbf{LUA}^{-1} = \mathbf{I}$, multiplicando ambos os lados da equação por \mathbf{L}^{-1} temos que

$$\mathbf{UA}^{-1} = \mathbf{L}^{-1}$$



Exemplo 5.8

Encontre a matriz inversa da seguinte matriz, por meio do método de eliminação de Gauss.

$$\mathbf{A} = \begin{bmatrix} 1 & 3 & 4 \\ 1 & 0 & 2 \\ -2 & 3 & 1 \end{bmatrix}$$

Etapa 01 - Escrever a matriz aumentada

$$\left(\begin{array}{ccc|ccc} 1 & 3 & 4 & 1 & 0 & 0 \\ 1 & 0 & 2 & 0 & 1 & 0 \\ -2 & 3 & 1 & 0 & 0 & 1 \end{array} \right)$$

Etapa 02 - Operações elementares sobre as linhas da matriz aumentada até obter uma matriz triangular superior e uma triangular inferior.

$$\left(\begin{array}{ccc|ccc} 1 & 3 & 4 & 1 & 0 & 0 \\ 0 & -3 & -2 & -1 & 1 & 0 \\ 0 & 0 & 3 & -1 & 3 & 1 \end{array} \right).$$

onde,

$$U = \begin{pmatrix} 1 & 3 & 4 \\ 0 & -3 & -2 \\ 0 & 0 & 3 \end{pmatrix}$$



Etapa 03 - Sabendo que $\mathbf{UA}^{-1} = \mathbf{L}^{-1}$, temos que

$$\mathbf{L}^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ -1 & 3 & 1 \end{pmatrix}.$$

implicando que,

$$\mathbf{A}^{-1} = \begin{pmatrix} 2/3 & -1 & -2/3 \\ 5/9 & -1 & -2/9 \\ -1/3 & 1 & 1/3 \end{pmatrix}.$$

Ou de maneira direta,

$$\mathbf{A} \cdot \mathbf{A}^{-1} = \mathbf{A}^{-1} \cdot \mathbf{A} = \mathbf{I}$$

Podemos concluir que a inversa pode ser calculada coluna a coluna, gerando soluções para vetores unitários.

$$\begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{bmatrix} \cdot \begin{bmatrix} d_1 \\ d_2 \\ d_3 \end{bmatrix}, \begin{bmatrix} d'_1 \\ d'_2 \\ d'_3 \end{bmatrix}, \begin{bmatrix} d''_1 \\ d''_2 \\ d''_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}, \begin{bmatrix} x'_1 \\ x'_2 \\ x'_3 \end{bmatrix}, \begin{bmatrix} x''_1 \\ x''_2 \\ x''_3 \end{bmatrix} = \begin{bmatrix} d_1 \\ d_2 \\ d_3 \end{bmatrix}, \begin{bmatrix} d'_1 \\ d'_2 \\ d'_3 \end{bmatrix}, \begin{bmatrix} d''_1 \\ d''_2 \\ d''_3 \end{bmatrix}$$

de modo que $\mathbf{X} = \mathbf{A}^{-1}$.



Exemplo 5.9

Um Algoritmo de criptografia (cifra de Hill) para textos utilizando a decomposição LU. Como exemplo, iremos criptografar o texto : *Resolver sistema*.

- **Passo 01:** Considere o seguinte conjunto, com $m = 29$, onde m representa a quantidade de letras e caracteres:

A = 1, B = 2, C = 3, D = 4, E = 5 F = 6 G = 7, H = 8, I = 9 J = 10, K = 11, L = 12, M = 13, N = 14, O = 15, P = 16, Q = 17, R = 18, S = 19, T = 20, U = 21, V = 22, W = 23, X = 24, Y = 25, Z = 26, Espaço = 27, # = 28 e @ = 0.

Para o texto escolhido temos que:

Resolver sistema = 18, 5, 19, 15, 12, 22, 5, 18, 27, 19, 9, 19, 20, 5, 13, 1
que escrito na forma matricial é equivalente a,

$$T = \begin{pmatrix} R & E & S & O \\ L & V & E & R \\ & S & I & S \\ T & E & M & A \end{pmatrix} = \begin{pmatrix} 18 & 5 & 19 & 15 \\ 12 & 22 & 5 & 18 \\ 27 & 19 & 9 & 19 \\ 20 & 5 & 13 & 1 \end{pmatrix}$$

Passo 02 : Escolha uma matriz referente a chave K da criptografia. No caso, iremos escolher a seguinte matriz,

$$K = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 2 & 4 & 8 & 16 \\ 3 & 9 & 27 & 81 \\ 4 & 16 & 64 & 256 \end{pmatrix}$$

de modo que $MDC(\det(K), m) = 1$. Esta condição garante a existência de uma matriz inversa $K^{-1}(\text{mod } m)$.

- **Passo 03:** Cifrar o texto por meio da seguinte operação matricial $C = (K T)(\text{mod } m)$,



$$C = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 2 & 4 & 8 & 16 \\ 3 & 9 & 27 & 81 \\ 4 & 16 & 64 & 256 \end{pmatrix} \begin{pmatrix} 18 & 5 & 19 & 15 \\ 12 & 22 & 5 & 18 \\ 27 & 19 & 9 & 19 \\ 20 & 5 & 13 & 1 \end{pmatrix} = \begin{pmatrix} 77 & 51 & 46 & 53 \\ 620 & 330 & 338 & 270 \\ 2511 & 1131 & 1398 & 801 \\ 7112 & 2868 & 4060 & 1820 \end{pmatrix} \text{mod } 29$$

$$C = \begin{cases} 77 = 2 \times 29 + 19 \\ 620 = 21 \times 29 + 11 \\ 2511 = 86 \times 29 + 17 \\ \vdots \\ 1820 = 62 \times 29 + 22 \end{cases} = \begin{pmatrix} 19 & 22 & 17 & 24 \\ 11 & 11 & 19 & 9 \\ 17 & 0 & 6 & 18 \\ 7 & 26 & 0 & 22 \end{pmatrix} = \begin{pmatrix} S & V & Q & X \\ K & K & S & I \\ Q & @ & F & R \\ G & Z & @ & V \end{pmatrix}$$

onde $SVQXKKSQ@FRGZ@V$ corresponde ao **texto cifrado**.

- **Passo 04** : Decompor a matriz $K = LU$.

$$K = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 3 & 3 & 1 & 0 \\ 4 & 6 & 4 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 2 & 6 & 14 \\ 0 & 0 & 6 & 36 \\ 0 & 0 & 0 & 24 \end{pmatrix}$$

onde $SVQXKKSQ@FRGZ@V$ corresponde ao **texto cifrado**.

- **Passo 05** : Criptografar o texto cifrado por meio de $X = (L^{-1}C) \pmod{m}$

$$X = \begin{pmatrix} 1 & 0 & 0 & 0 \\ -2 & 1 & 0 & 0 \\ 3 & -3 & 1 & 0 \\ -4 & 6 & -4 & 1 \end{pmatrix} \begin{pmatrix} 19 & 22 & 17 & 24 \\ 11 & 11 & 19 & 9 \\ 17 & 0 & 6 & 18 \\ 7 & 26 & 0 & 22 \end{pmatrix} = \begin{pmatrix} 19 & 22 & 17 & 24 \\ -27 & -33 & -15 & -39 \\ 41 & 33 & 0 & 63 \\ -71 & 4 & 22 & -99 \end{pmatrix} \text{mod } 29$$

$$= \begin{pmatrix} 19 & 22 & 17 & 24 \\ 2 & 25 & 14 & 19 \\ 12 & 4 & 0 & 5 \\ 16 & 4 & 22 & 24 \end{pmatrix} = \begin{pmatrix} S & V & Q & X \\ B & Y & N & S \\ L & D & @ & E \\ P & D & V & X \end{pmatrix} \Rightarrow SVQXBYNSLD@EPDVX.$$



Portanto, $SVQXBYNSLD@EPDV$ X corresponde efetivamente ao **texto criptografado**.

Antes de explicar o processo de decriptografia, se faz necessário o entendimento da inversa de uma matriz modular.

Exemplo 5.10

- Calcular o determinante da matriz K :

$$\det(K) = \prod_{i=1}^n u_{ii}$$

e $MDC(\det(K), m) = 1$ de modo a garantir que a matriz U é invertível módulo m . No exemplo anterior temos que $\det(K) = 288$ e $MDC(288, 29) = 1$. Portanto, a matriz U é invertível módulo 29.

– Encontrar $(\det(K))^{-1}$, onde

$$(\det(K))(\det(K))^{-1} \equiv 1 \pmod{m}.$$

Como temos que $\det(K) = 288 \rightarrow 288.x \equiv 1 \pmod{29}$, onde x corresponde ao inverso multiplicativo modular de 288, então

$$\left\{ \begin{array}{l} x = 1 \rightarrow 288 = 1 \times 29 + 259 \\ x = 2 \rightarrow 288.2 = 976 = 33 \times 29 + 19 \\ x = 3 \rightarrow 288.3 = 864 = 29 \times 29 + 23 \\ x = 4 \rightarrow 288.4 = 1152 = 39 \times 29 + 21 \\ \vdots \\ x = 14 \rightarrow 288.14 = 4032 = 139 \times 29 + 1. \end{array} \right.$$

Portanto, $(\det(K))^{-1} = 14$.

– Sabemos que $U^{-1} = (\det(K))^{-1} \cdot adj(U) \pmod{m}$, onde $adj(U)$ corresponde a adjunta da matriz U . Ou seja,



$$\begin{aligned}
 U^{-1} &= 14 \begin{pmatrix} 288 & -144 & 96 & -72 \\ 0 & 144 & -144 & 132 \\ 0 & 0 & 48 & -72 \\ 0 & 0 & 0 & 12 \end{pmatrix} \\
 &= \begin{pmatrix} 4032 & -2016 & 1344 & -1008 \\ 0 & 2016 & -2016 & 1848 \\ 0 & 0 & 672 & -1008 \\ 0 & 0 & 0 & 168 \end{pmatrix}_{\text{mod } 29}
 \end{aligned}$$

de modo que,

$$U^{-1} = \begin{pmatrix} 1 & 14 & 10 & 7 \\ 0 & 15 & 14 & 21 \\ 0 & 0 & 5 & 7 \\ 0 & 0 & 0 & 23 \end{pmatrix}$$

- **Passo 06** : Obtenção do texto original por meio da seguinte operação,

$$T_o = (U^{-1}X)(\text{mod } m).$$

$$T_o = \begin{pmatrix} 1 & 14 & 10 & 7 \\ 0 & 15 & 14 & 21 \\ 0 & 0 & 5 & 7 \\ 0 & 0 & 0 & 23 \end{pmatrix} \begin{pmatrix} 19 & 22 & 17 & 24 \\ 2 & 25 & 14 & 19 \\ 12 & 4 & 0 & 5 \\ 16 & 4 & 22 & 24 \end{pmatrix}$$



$$\begin{aligned}
&= \begin{pmatrix} 279 & 440 & 367 & 508 \\ 534 & 515 & 672 & 859 \\ 172 & 48 & 154 & 193 \\ 368 & 92 & 506 & 552 \end{pmatrix} \equiv_{\text{mod } 29} \begin{pmatrix} 18 & 5 & 19 & 15 \\ 12 & 22 & 5 & 18 \\ 27 & 19 & 9 & 19 \\ 20 & 5 & 13 & 1 \end{pmatrix} \\
&= \begin{pmatrix} R & E & S & O \\ L & V & E & R \\ & S & I & S \\ T & E & M & A \end{pmatrix} \implies \text{RESOLVER SISTEMA.}
\end{aligned}$$

Que corresponde ao **texto original**.

Se métodos diretos para sistemas lineares, que encontram a resposta exata (dentro do erro de arredondamento) em um número finito de etapas, estão disponíveis, então por que desenvolver métodos iterativos? Muitas vezes a matriz de coeficientes do sistema é uma matriz completamente arbitrária, em alguns casos esparsa. Assim, o sistema pode ser resolvida de maneira mais eficiente via métodos iterativos. Existem inúmeras aplicações em que tais métodos são rápidos e eficientes computacionalmente.

MÉTODOS ITERATIVOS

Passamos agora a uma descrição de métodos iterativos para resolver sistemas de equações lineares. Esses métodos fornecem uma solução para um sistema por meio do limite de uma sequência vetorial, construída utilizando um processo denominado iterativo. Na literatura, um grande número de métodos iterativos são descritos com base em princípios diferentes. Como regra, os algoritmos computacionais utilizados na resolução de tais métodos são simples e convenientes computacionalmente. No entanto, podem surgir



casos em que, dependendo da aplicação a convergência pode ser lenta, tornando-o inviável na prática. Nos métodos iterativos, as equações são reorganizadas de forma a permitir operações iterativas até que a convergência seja alcançada. Se faz necessário uma estimativa inicial. Um método para calcular a solução única de um sistema $\mathbf{Ax} = \mathbf{b}$, $\mathbf{A} = (a_{ij})$, $i, j = 1, \dots, n$ e $\det(\mathbf{A}) \neq 0$ é denominado iterativo quando fornece uma sequência de soluções aproximadas a partir de iterações anteriores. Ou seja, a sequência de vetores $\mathbf{x}^{(1)}$, $\mathbf{x}^{(2)}$, \dots , $\mathbf{x}^{(k)}$, $\mathbf{x}^{(k+1)}$ é construída a partir da seguinte recorrência:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^k + \mathbf{H}^{(k+1)}(\mathbf{b} - \mathbf{Ax}^k), \quad k = 0, 1, \dots$$

onde $\mathbf{H}^{(1)}$, $\mathbf{H}^{(2)}$, \dots é uma certa sequência matricial e $\mathbf{x}^{(0)}$ a condição inicial, geralmente arbitrária. Uma escolha diferente da sequência matricial $\mathbf{H}^{(k)}$ representa diferentes processos iterativos. No caso, iremos considerar a seguinte expressão,

$$\mathbf{x}^{(k+1)} = \mathbf{Cx}^k + \mathbf{d}, \quad k = 0, 1, \dots$$

em que $\mathbf{C}_{n \times n}$ é a matriz das iterações e $\mathbf{d}_{n \times 1}$ é um vetor. De modo que,

$$\lim_{k \rightarrow \infty} \mathbf{x}^k = \bar{\mathbf{x}}$$

Teorema 5.1

Para um processo de aproximação sucessiva convergir é necessário e suficiente que

$$\|\mathbf{C}\| < 1.$$



MÉTODO ITERATIVO DE JACOBI - RICHARDSON

seja o SEL:

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1 \\ a_{12}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2 \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n = b_n \end{cases}$$

$$\begin{cases} x_1 = \frac{1}{a_{11}}(b_1 - (a_{12}x_2 + a_{13}x_3 + \cdots + a_{1n}x_n)) \\ x_2 = \frac{1}{a_{22}}(b_2 - (a_{21}x_1 + a_{23}x_3 + \cdots + a_{2n}x_n)) \\ \vdots \\ x_n = \frac{1}{a_{nn}}(b_n - (a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn-1}x_{n-1})) \end{cases}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} 0 & \frac{-a_{12}}{a_{11}} & \frac{-a_{13}}{a_{11}} & \cdots & \frac{-a_{1n}}{a_{11}} \\ \frac{-a_{22}}{a_{22}} & 0 & \frac{-a_{23}}{a_{22}} & \cdots & \frac{-a_{2n}}{a_{22}} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \frac{-a_{n1}}{a_{nn}} & \frac{-a_{n2}}{a_{nn}} & \cdots & \cdots & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} \frac{b_1}{a_{11}} \\ \frac{b_2}{a_{22}} \\ \vdots \\ \frac{b_n}{a_{nn}} \end{bmatrix}$$

Assim, podemos escrever o sistema como sendo $\mathbf{x} = \mathbf{C}\mathbf{x} + \mathbf{d}$, onde:

$$c_{ij} = \begin{cases} 0 & \text{se } i = j \\ \frac{-a_{ij}}{a_{ii}} & \text{se } i \neq j \end{cases} \quad i, j = 1, 2, \dots, n$$

e

$$d_{ij} = \frac{b_i}{a_{ii}} \quad i, j = 1, 2, \dots, n$$

Dessa forma podemos escrever o método iterativo de Jacobi - Richardson:

$$\mathbf{x}^{(k+1)} = \mathbf{C}\mathbf{x}^{(k)} + \mathbf{d} \quad k = 0, 1, \dots$$



ou seja,

...

$$\begin{cases} x_1^{(k+1)} = \frac{1}{a_{11}} \left(b_1 - (a_{12}x_2^{(k)} + a_{13}x_3^{(k)} + \dots + a_{1n}x_n^{(k)}) \right) \\ x_2^{(k+1)} = \frac{1}{a_{22}} \left(b_2 - (a_{21}x_1^{(k)} + a_{23}x_3^{(k)} + \dots + a_{2n}x_n^{(k)}) \right) \\ \vdots \\ x_n^{(k+1)} = \frac{1}{a_{nn}} \left(b_n - (a_{n1}x_1^{(k)} + a_{n2}x_2^{(k)} + \dots + a_{nn-1}x_{n-1}^{(k)}) \right) \end{cases}$$

Que na forma geral, pode ser escrito como:

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left[b_i - \left(\sum_{j=1; j \neq i}^n a_{ij}x_j^{(k)} \right) \right], \quad i = 1, 2, \dots, n$$

CONVERGÊNCIA

Sejam uma norma matricial consistente e $\mathbf{x}^{(0)} \in \mathbb{R}^n$ a solução inicial. Se $\|\mathbf{C}\| < 1$, então a sequência de soluções será tal que $\mathbf{x}^{(k+1)} = \mathbf{C}\mathbf{x}^{(k)} + \mathbf{d}$, $k = 0, 1, \dots$ converge para a solução \mathbf{x}^* do sistema $\mathbf{A}\mathbf{x} = \mathbf{b}$.

Uma possível condição de parada:

$$E_r = \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|_{\infty}}{\|\mathbf{x}^{(k+1)}\|_{\infty}} \leq \varepsilon$$

ALGORITMO - MÉTODO DE JACOBI

Dados : $A_{n \times n}$, $b_{n \times 1}$, $x^{(0)}$, max , tol

Para $k = 1 : max$, faça

Para $i = 1 : n$, faça



$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1, j \neq i}^n a_{ij} x_j^{(k)} \right)$$

Se $Er \leq tol$, então

$$\mathbf{x} = \mathbf{x}^{(k+1)}$$

Senão, se $k = max$, então

não houve convergência

Exemplo 5.11

Usando o método de Jacobi - Richardson, determine uma solução aproximada para o seguinte SEL, com aproximação inicial $\mathbf{x}^{(0)} = [0 \ 0 \ 0]^T$ e precisão de $\varepsilon = 0, 01$.

$$\begin{cases} 10x_1 + 2x_2 + x_3 = 14 \\ x_1 + 5x_2 + x_3 = 11 \\ 2x_1 + 3x_2 + 10x_3 = 8 \end{cases}$$

Exemplo 5.12

Usando o método de Jacobi - Richardson, determine uma solução aproximada para o seguinte SEL, com aproximação inicial $\mathbf{x}^{(0)} = [0, 7 - 1, 6 \ 0, 6]^T$ e precisão de $\varepsilon = 10^{-2}$.

$$\begin{cases} x_1 + 5x_2 + x_3 = -8 \\ 10x_1 + 2x_2 + x_3 = 7 \\ 2x_1 + 3x_2 + 10x_3 = 6 \end{cases}$$



MÉTODO DE GAUSS - SEIDEL

Esse método difere do processo de Jacobi por utilizar para o cálculo de uma componente de $x^{(k+1)}$ o valor mais recente das demais componentes.

Forma geral:

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left[b_i - \left(\sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} + \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right) \right], \quad i = 1, 2, \dots, n$$

TESTES DE CONVERGÊNCIA

O método converge se:

- O critério das linhas for satisfeito:

$$\|\mathbf{A}\|_1 = \max_{1 \leq i \leq n} \sum_{j=1, j \neq i}^n |a_{ij}^*| < 1$$

- Ou o critério das colunas for satisfeito:

$$\|\mathbf{A}\|_\infty = \max_{1 \leq j \leq n} \sum_{i=1, i \neq j}^n |a_{ij}^*| < 1$$

onde

$$a_{ij}^* = \frac{a_{ij}}{a_{ii}}$$

- Ou a matriz é **Extritamente** Diagonalmente predominante:

$$\sum_{j=1, j \neq i}^n |a_{ij}| < |a_{ii}|, \quad i = 1, 2, \dots, n$$



- Ou atender ao critério de Sassenfeld:

$$\beta = \max_{1 \leq i \leq n} \beta_i < 1$$

Critério de Sassenfeld:

$$\beta = \max_{1 \leq i \leq n} \beta_i < 1$$

onde

$$\beta_1 = \sum_{j=2}^n |a_{1j}^*|$$

e

$$\beta_i = \sum_{j=1}^{i-1} |a_{ij}^*| \beta_j + \sum_{j=i+1}^n |a_{ij}^*|, \quad i = 2, \dots, n.$$

Se $\beta < 1$, a sequência $x^{(k)}$, gerada pelo método iterativo de Gauss - Seidel, converge para o solução do sistema.

ALGORITMO - MÉTODO DE GAUSS - SEIDEL

Dados : $A_{n \times n}$, $b_{n \times 1}$, $x^{(0)}$, max

Para $k = 1 : max$, faça

Para $i = 1 : n$, faça

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right)$$



Exemplo 5.13

$$\begin{cases} 5x_1 + x_2 + x_3 = 5 \\ 3x_1 + 4x_2 + x_3 = 6 \\ 3x_1 + 3x_2 + 6x_3 = 0 \end{cases}$$

MÉTODOS DIRETOS X MÉTODOS ITERATIVOS

MÉTODOS DIRETOS

Recomendados para sistemas de pequeno porte com matrizes de coeficiente densas.

MÉTODOS ITERATIVOS

Bastante vatajosos para sistemas de grande porte cuja matriz de coeficiente seja esparsa. No entanto, é necessário verificar as condições de convergência.

EXERCÍCIOS PROPOSTOS

Exercício 5.1

Considere o SEL $A.x = b$,

$$\begin{bmatrix} 1 & \alpha & 3 \\ \alpha & 1 & 4 \\ 5 & 2 & 1 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -2 \\ -3 \\ 4 \end{bmatrix}$$



- Para que valores de α , a matriz A pode ser decomposta no produto LU ?
- Considere $\alpha = 1$ e resolva o sistema por Eliminação de Gauss.

Exercício 5.2

Resolva o SEL adequadamente, utilizando decomposição de *Cholesky*

$$\begin{cases} 2x_1 + yx_2 - x_3 = 3 \\ x_1 + 10x_2 + zx_3 = 6 \\ kx_1 + 2x_2 + 4x_3 = -6 \end{cases}$$

Exercício 5.3

Num determinado circuito elétrico, as correntes i_1 , i_2 e i_3 passam através das impedâncias z_1 , z_2 e z_3 por meio do seguinte SEL:

$$\begin{cases} i_1 + i_2 + i_3 = v_1 - v_3 \\ z_1 i_1 - z_2 i_2 = v_1 - v_2 \\ z_2 i_1 - z_3 i_3 = v_2 - v_3 \end{cases}$$

com, $v_1 + v_2 + v_3 = 122V$, $v_1 - v_3 = 0$, $v_1 - v_2 = 70V$, $z_1 = 6\Omega$, $z_2 = 8\Omega$ e $z_3 = 4\Omega$.

Determine as correntes no circuito.



Exercício 5.4

Um estudante de Engenharia deseja montar um computador, no entanto ainda faltam 3 componentes. A, B e C. Se adquirir, respectivamente:

- 4, 5 e 6 componentes, gastará 1.700, 00;
- 5, 2 e 10 componentes, gastará 2.120, 00;
- 6, 6 e 4 componentes, gastará 1.670, 00.

Determine o preço de cada componente.

Exercício 5.5

Resolva o SEL pelo método iterativo de Gauss - Seidel, de modo que o erro relativo sejam menor que 10^{-2} ,

$$\begin{cases} x_1 + x_2 - 5x_3 = 4 \\ x_1 - 10x_2 - x_3 = 2 \\ 5x_1 + x_2 + x_3 = 1 \end{cases}$$

Exercício 5.6

Ordene o seguinte SEL de modo que o critério de convergência de Sassenfeld seja atendido.

$$\begin{cases} 3x_1 + 3x_2 - 5x_3 = 2 \\ 10x_1 + 3x_2 + 2x_3 = -20 \\ 5x_1 + 5x_2 - 3x_3 = 10 \end{cases}$$



Exercício 5.7

Sistemas idealizados de massa-mola desempenham um papel importante na mecânica e em outros problemas da engenharia. A Figura abaixo mostra tal sistema. Assim que são soltas, as massas são puxadas para baixo pela força da gravidade. Observe que o deslocamento de cada mola é medido em coordenadas locais relativas à sua posição inicial.

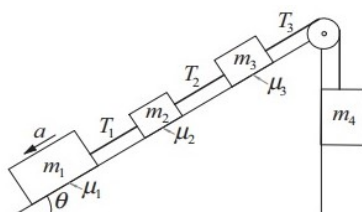
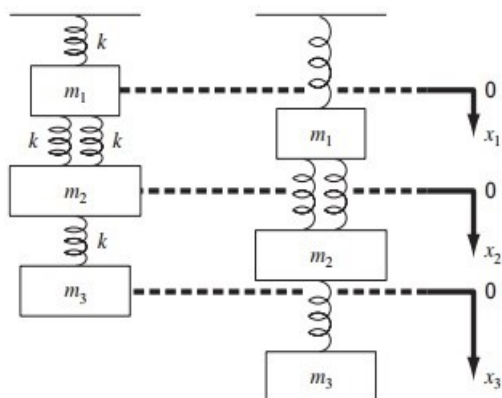
- Usando a segunda lei de Newton e a lei de Hooke modele um sistema de equações lineares;
- considerando as derivadas nulas e, $m_1 = 2\text{kg}$, $m_2 = 3\text{kg}$, $m_3 = 2,5\text{kg}$ e $k = 10^4\text{kg}$, use a decomposição LU para determinar os deslocamentos.

Exercício 5.8

Quatro blocos de diferentes massas m_i são conectados por cordas de massa desprezível. Três blocos estão em um plano inclinado, os coeficientes de atrito entre os blocos e plano são μ_i . As equações de movimento para o bloco são definidas por,

$$\begin{cases} T_1 + m_1 a = m_1 g (\sin\theta - \mu_1 \cos\theta) \\ T_2 - T_1 + m_2 a = m_2 g (\sin\theta - \mu_2 \cos\theta) \\ T_3 - T_2 + m_3 a = m_3 g (\sin\theta - \mu_3 \cos\theta) \\ m_4 a - T_3 = -m_4 g \end{cases}$$





onde T_i representa o torque e a a aceleração do sistema. Determine a e T_i para um ângulo de $\theta = \pi/4$, $g \approx 9,8 \text{ m/s}^2$, $\mu = [0, 25 \ 0, 3 \ 0, 2]^T$ e $m = [10 \ 4 \ 5 \ 6]^T$.



AJUSTE DE CURVAS

Um problema básico na engenharia é ajustar um modelo matemático a um conjunto de dados sujeito a erros (possivelmente erros de medição e no modelo). Ou seja, dados n pares ordenados $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ devem ser ajustados a um modelo descrito por uma função $f(x)$ (ou curva) tal que,

$$\begin{cases} y_1 \approx f(x_1), \\ y_2 \approx f(x_2), \\ \vdots \\ y_n \approx f(x_n) \end{cases}$$

represente a tendência geral dos dados. O tipo de função (polinomial, senoidal, exponencial, de base radial, etc) é determinada pela natureza do problema. Claramente, quanto mais dados (observações) estiverem disponíveis, mais preciso será o modelo proposto no ajuste do modelo. Na engenharia um problema de ajuste de curvas ou regressão, pode ser visto por exemplo, como um meio de filtrar um sinal atenuando ou até eliminar o ruído (muitas vezes, indesejável) presente ao sinal. Questões como essas, estão diretamente relacionadas à Estatística e a teoria de aproximação de funções. Na formulação da teoria, os modelos de regressão são frequentemente ajustados utilizando uma abordagem denominada de mínimos quadrados. Ou seja, o objetivo principal é obter uma função que minimize o erro médio quadrático entre os dados. Modelos de regressão linear são algoritmos supervisionados, em

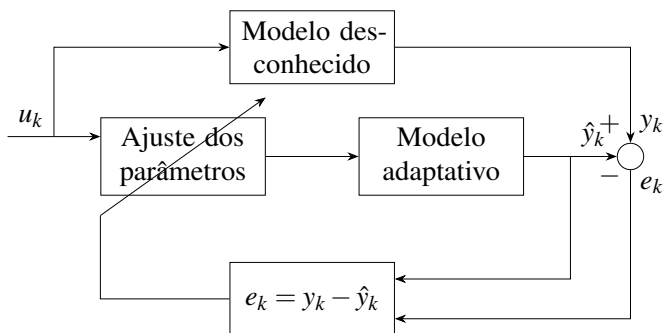


que existe uma relação entre a variável objetivo e as variáveis independentes do sistema. O modelo geralmente aprende os padrões dos dados apresentados e fornece uma equação matemática que mostra a tendência ou comportamento do sistema desconhecido. Estes tipo de modelo são frequentemente utilizados na resolução de dois problemas clássicos: Classificação e Regressão.

REGRESSÃO POR MÍNIMOS QUADRADOS

Seja o modelo apresentado na seguinte figura,

Figura 6.1 – Clássico problema de regressão

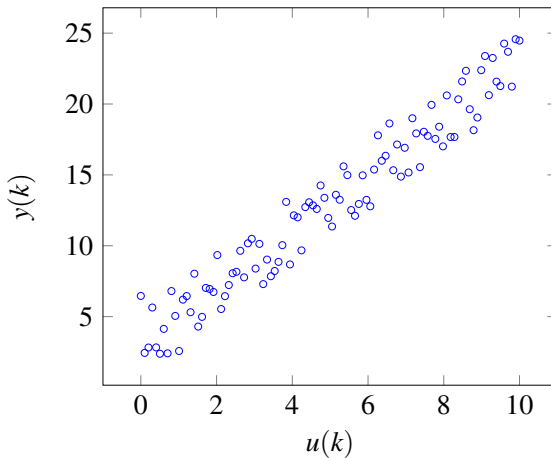


onde um sinal u_k é apresentado a um sistema desconhecido e ao mesmo tempo a um modelo matemático adaptativo. De modo a obtermos dois sinais, y_k fornecido pelo sistema desconhecido e \hat{y}_k fornecido pelo modelo adaptativo. Ao compararmos os respectivos sinais, obtemos um erro $e_k = y_k - \hat{y}_k$. O objetivo num problema dessa natureza é ajustar a saída do modelo adaptativo \hat{y}_k de modo ao mesmo ficar o mais próximo possível do sinal apresentado pelo sistema desconhecido y_k , minimizando o erro quadrático. Uma vez atingido o objetivo, dizemos que o modelo matemático adaptativo, ajustou-se ao sistema desconhecido, representando assim toda a



dinâmica do mesmo. Ou seja, o modelo “identifica” o sistema. Suponha que ao ser apresentado um sinal u_k ao sistema desconhecido da figura (6.1), foi fornecido o seguinte sinal de saída y_k .

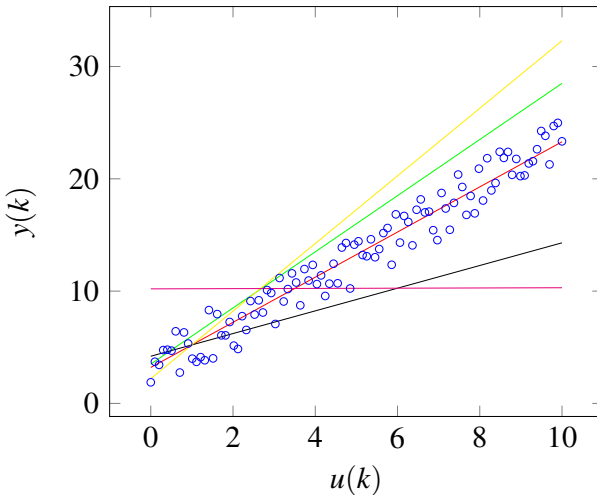
Figura 6.2 – Exemplo aplicado ao sistema da figura (6.1)



Uma pergunta natural é, dentre os modelos matemático apresentados na figura (6.2), qual o que “melhor” se ajusta ao sistema? Baseado em qual métrica? Como escolher o “melhor” modelo? Por exemplo, ao analisarmos a seguinte figura, o que podemos concluir?



Figura 6.3 – Diferentes modelos de ajuste aos dados da figura (6.2)



Pela tendência dos dados, claramente vemos que uma hipótese razoável seria um modelo matemático, tal como:

$$\hat{y}(u_k) = f(w_k, u_k) = w_0 + w_1 u_k$$

com isso, nosso problema se resume a encontrar os valores w_0 e w_1 que melhor ajustam o modelo. Mas, ainda resta uma dúvida. qual a métrica utilizada para medir a qualidade do ajuste?. Ou seja, dado o conjunto de dados de entrada e saída do modelo (desconhecido), como identificá-lo ou representá-lo pelo método de regressão? Para isso iremos utilizar um **funcional de custo**, definido como:

$$J = \frac{1}{2} \sum_{k=1}^m (e_k)^2$$

onde $e_k = y_k - \hat{y}_k$ representa o erro do processo. De modo que para o problema proposto, podemos reescrever a expressão (6.1) como,



$$J = \frac{1}{2} \sum_{k=1}^m (y_k - \hat{y}(u_k))^2 \Rightarrow J = \frac{1}{2} \sum_{k=1}^m (y_k - (w_0 + w_1 u_k))^2.$$

Para minimizar o erro precisamos encontrar os valores de w_0 e w_1 que nos dão o ponto de mínimo na superfície quadrática dada pelo funcional de custo. Para isso, utilizaremos o conceito do vetor gradiente, iremos operar da seguinte forma:

$$\frac{\partial J}{\partial w_0} = 0 \text{ e } \frac{\partial J}{\partial w_1} = 0$$

obtendo as seguintes expressões,

$$\frac{\partial J}{\partial w_0} = \left[\sum_{k=1}^m (y_i - (w_0 + w_1 u_k)) \right] (-1) = 0$$

$$\frac{\partial J}{\partial w_1} = \left[\sum_{k=1}^m (y_i - (w_0 + w_1 u_k)) \right] \left(- \sum_{k=1}^m x_k \right) = 0$$

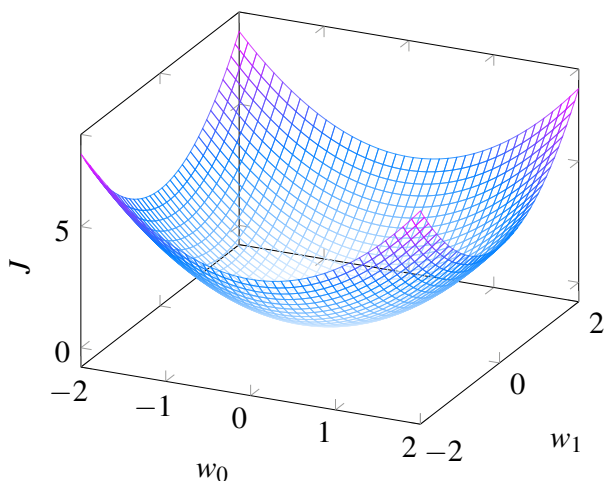
que corresponde a um sistema de equações lineares, que escrito na forma matricial, corresponde a

$$\begin{pmatrix} m & \sum_{k=1}^m u_k \\ \sum_{k=1}^m u_k & \sum_{k=1}^m u_k^2 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} = \begin{pmatrix} \sum_{k=1}^m y_k \\ \sum_{k=1}^m y_k \cdot u_k \end{pmatrix} \rightarrow \begin{pmatrix} \langle \mathbf{1}, \mathbf{1} \rangle & \langle \mathbf{1}, \mathbf{u} \rangle \\ \langle \mathbf{u}, \mathbf{1} \rangle & \langle \mathbf{u}, \mathbf{u} \rangle \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} = \begin{pmatrix} \langle \mathbf{1}, \mathbf{y} \rangle \\ \langle \mathbf{y}, \mathbf{u} \rangle \end{pmatrix} \quad (6.2)$$

em que $\mathbf{1} = (1_1, 1_2, \dots, 1_m)^T$, $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$ e $\mathbf{u} = (u_1, u_2, \dots, u_m)^T$ pertencem ao \mathbb{R}^m e $\mathbf{w} = (w_0, w_1)^T \in \mathbb{R}^2$.



Figura 6.4 – Gráfico do funcional de custo para duas variáveis w_0 e w_1



Exemplo 6.1

Obter a reta que melhor se ajusta os dados da seguinte tabela

Tabela 6.1 – Dados do exemplo 6.1.

0	1	2	3	4
0,98	-3,01	-6,99	-11,01	-15

$$\begin{pmatrix} 5 & 10 \\ 10 & 30 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} = \begin{pmatrix} -35,03 \\ -110,02 \end{pmatrix} \rightarrow \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} = \begin{pmatrix} 0,986 \\ -3,996 \end{pmatrix}$$

$$\hat{y} = 0,986 \cdot \mathbf{1} - 3,996 \cdot \mathbf{u}$$

Exemplo 6.2

Obter o modelo que melhor se adapta ao seguinte conjunto de dados



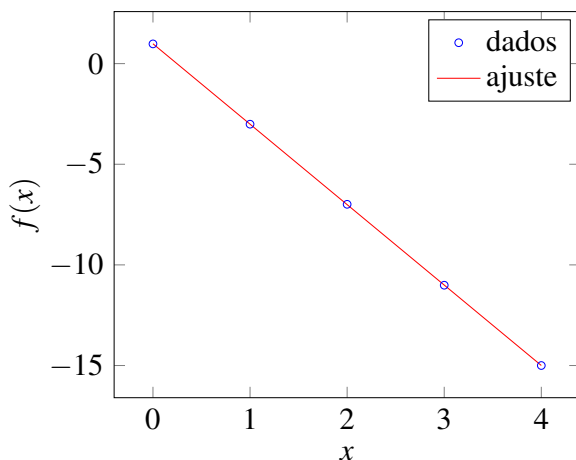
Tabela 6.2 – Dados do exemplo 6.2.

x_k	2	3	4	5	6	7	8
y_k	2,85	3,53	4,30	5,14	6,03	6,96	7,93

$$\begin{pmatrix} 7 & 35 \\ 35 & 203 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} = \begin{pmatrix} 36,74 \\ 207,53 \end{pmatrix} \rightarrow \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} = \begin{pmatrix} 0,9932 \\ 0,8510 \end{pmatrix}$$

$$\hat{y}_k = w_0 + w_1 \cdot x_k = 0,9932 + 0,8510 \cdot x_k, \quad k = 0, 1, \dots, 6$$

Figura 6.5 – Gráfico do exemplo 6.1.



Exemplo 6.3

Uma pesquisa de mercado entre os anos 1992 e 2015 mostrou a evolução na valorização (em milhares de dólares) de uma certa classe de imóveis em Los Angeles/CA. A partir dos dados apresentados, obtenha um modelo preditivo para o mercado imobiliário da cidade no período.



i	Data	US\$	Data	US\$	Data	US\$
1	1992/01/01	92,45	2000/01/01	100,00	2008/03/01	207,11
2	1992/06/01	90,41	2000/06/01	106,36	2008/07/01	192,55
3	1993/01/01	84,88	2001/02/01	111,32	2009/01/01	166,55
4	1993/07/01	81,11	2001/07/01	117,36	2009/06/01	160,90
5	1994/02/01	76,86	2002/01/01	121,45	2010/01/01	172,97
6	1994/07/01	76,95	2002/06/01	131,59	2010/07/01	176,27
7	1995/01/01	75,91	2003/01/01	144,27	2011/02/01	168,23
8	1995/08/01	75,11	2003/05/01	152,27	2011/06/01	169,66
9	1995/12/01	73,60	2004/02/01	180,49	2012/01/01	160,76
10	1996/08/01	74,9	2004/06/01	206,38	2012/08/01	173,01
11	1997/01/01	73,91	2005/03/01	226,75	2013/01/01	180,23
12	1997/05/01	75,72	2005/07/01	246,37	2013/06/01	202,08
13	1998/01/01	80,28	2006/02/01	267,75	2014/02/01	215,25
14	1998/06/01	86,97	2006/06/01	273,22	2014/07/01	224,56
15	1999/02/01	91,46	2007/01/01	268,68	2015/01/01	225,94
16	1999/08/01	98,30	2007/06/01	262,12	2015/06/01	237,14

CASO GERAL POLINOMIAL

Se o modelo que desejamos ajustar for uma combinação linear de várias funções $p_i(x)$ na forma

$$y \approx f(x) = \sum_{i=0}^n w_i p_i(x) = w_0 p_0(x) + w_1 p_1(x) + \dots + w_n p_n(x)$$

onde as constantes w_i 's são desconhecidas, e os p_i 's são as denominadas funções de base do tipo polinomial.



Exemplo 6.4

Seja $p_i(x) = x^i$, $i = 1, \dots, n$ uma classe de funções que constituem uma base $\{1, x, x^2, \dots, x^n\}$ no espaço n dimensional de funções polinomiais. De modo que qualquer função neste espaço pode ser escrito como uma combinação linear de funções da base polinomial. Ou seja,

$$f(x) = \sum_{i=0}^n w_i x^i = w_0 + w_1 x + w_2 x^2 + \dots + w_n x^n$$

Para resolver um problema com funções de base polinomiais, em que dado o conjunto de pares ordenados

$$(x_i; f(x_i))_{i=0}^m,$$

obter os valores dos parâmetros desconhecidos, se resume a resolver o seguinte sistema de equações lineares,

$$\begin{cases} w_0 + w_1 x_0 + w_2 x_0^2 + \dots + w_n x_0^n = f(x_0) \\ w_0 + w_1 x_1 + w_2 x_1^2 + \dots + w_n x_1^n = f(x_1) \\ w_0 + w_1 x_2 + w_2 x_2^2 + \dots + w_n x_2^n = f(x_2) \\ \vdots \\ w_0 + w_1 x_m + w_2 x_m^2 + \dots + w_n x_m^n = f(x_m) \end{cases}$$

que escrito na forma matricial equivale a,

$$\begin{pmatrix} 1 & x_0 & x_0^2 & \dots & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^n \\ 1 & x_2 & x_2^2 & \dots & x_2^n \\ \vdots & \vdots & \dots & \ddots & \vdots \\ 1 & x_m & x_m^2 & \dots & x_m^n \end{pmatrix}_{m \times n} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix}_{n \times 1} = \begin{pmatrix} f(x_0) \\ f(x_1) \\ f(x_2) \\ \vdots \\ f(x_m) \end{pmatrix}_{m \times 1} \Rightarrow \mathbf{X} \cdot \mathbf{w} = \mathbf{f}.$$



Utilizando o funcional de custo,

$$J = \frac{1}{2} \|\mathbf{f} - \mathbf{X}\mathbf{w}\|^2 = \frac{1}{2} (\mathbf{f} - \mathbf{X}\mathbf{w})^T (\mathbf{f} - \mathbf{X}\mathbf{w}) = \frac{1}{2} (\mathbf{f}^T \mathbf{f} - \mathbf{f}^T \mathbf{X}\mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{f} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w})$$

fazendo,

$$\frac{\partial J}{\partial \mathbf{w}} = \mathbf{0}$$

obteremos:

$$-\mathbf{X}^T \mathbf{f} + \mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{0}$$

que representa as **equações normais** que fornecem o vetor de parâmetros w 's que minimizam o erro quadrático. Assim,

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{f}$$

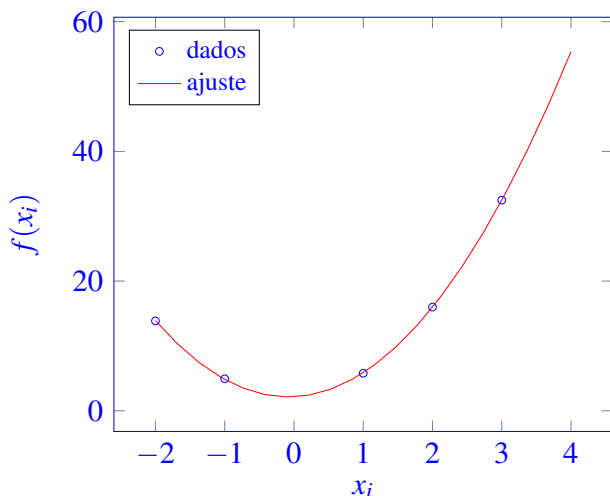
Exemplo 6.5

Obtenha uma aproximação polinomial do tipo,

$$f(x) = w_0 + w_1x + w_2x^2$$



Figura 6.6 – Gráfico do exemplo 6.6.



para o conjunto de dados da seguinte tabela

$$\begin{pmatrix} 1 & -2 & 4 \\ 1 & -1 & 1 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} 13,86 \\ 04,93 \\ 05,79 \\ 15,99 \\ 32,48 \end{pmatrix}$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{f}$$

$$\begin{pmatrix} 5 & 3 & 19 \\ 3 & 19 & 27 \\ 19 & 27 & 115 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} 3,05 \\ 102,56 \\ 422,44 \end{pmatrix}$$



$$\rightarrow \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} 2,1541 \\ 0,5154 \\ 3,1965 \end{pmatrix}$$

$$f(x) = 2,1541 + 0,5154x + 3,1965x^2.$$

Para medir a qualidade do ajuste iremos utilizar como referência o vetor erro residual, ou seja

$$\mathbf{r} = \mathbf{f} - \mathbf{Aw}$$

Iremos adotar 2 (duas) formas para expressar um quantitativo para o residual. A norma euclidiana,

$$\|\mathbf{r}\|_2 = \sqrt{r_1^2 + r_2^2 + \dots + r_n^2}$$

e, a raiz da média da norma do erro quadrático.

$$RMSE = \sqrt{m} \frac{\|\mathbf{r}\|_2}{m} \quad (6.6)$$

Para o exemplo (6.3) temos,

$$\mathbf{r} = \begin{pmatrix} -0,0492 \\ 0,0948 \\ -0,0759 \\ 0,01915 \\ 0,01129 \end{pmatrix} \rightarrow \|\mathbf{r}\|_2 = 0,1329 \text{ e } RMSE = 0,0594.$$

FUNÇÕES DE BASE

As funções polinomiais consideradas na seção anterior é um exemplo particular de um modelo em que a variável de entrada



x , e as funções de base assumem a forma de potências de x de modo que $g_k(x) = x^k$. Um limitação das funções de base polinomial é que as mesmas são funções globais com relação a variável de entrada, de modo que mudanças em alguma região do espaço de entrada, afetam todas as demais regiões. Assim, se o modelo que desejamos ajustar for uma combinação linear de várias funções de base $g_i(x)$ na forma

$$y \approx f(x) = \sum_{i=0}^n w_i g_i(x) = w_0 g_0(x) + w_1 g_1(x) + \cdots + w_n g_n(x)$$

onde as constantes w 's são desconhecidas, e as g 's são as funções de base em um espaço de funções. Que pode ser escrito como o seguinte sistema de equações lineares,

$$\begin{cases} w_0 g_0(x_0) + w_1 g_1(x_0) + w_2 g_2(x_0) + \cdots + w_n g_n(x_0) = f(x_0) \\ w_0 g_0(x_1) + w_1 g_1(x_1) + w_2 g_2(x_1) + \cdots + w_n g_n(x_1) = f(x_1) \\ w_0 g_0(x_2) + w_1 g_1(x_2) + w_2 g_2(x_2) + \cdots + w_n g_n(x_2) = f(x_2) \\ \vdots \\ w_0 g_0(x_n) + w_1 g_1(x_n) + w_2 g_2(x_n) + \cdots + w_n g_n(x_n) = f(x_n) \end{cases}$$

na forma matricial,

$$\begin{pmatrix} g_0(x_0) & g_1(x_0) & g_2(x_0) & \cdots & g_n(x_0) \\ g_0(x_1) & g_1(x_1) & g_2(x_1) & \cdots & g_n(x_1) \\ \vdots & \vdots & \ddots & \cdots & \vdots \\ g_0(x_n) & g_1(x_n) & g_2(x_n) & \cdots & g_n(x_n) \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{pmatrix} = \begin{pmatrix} f(x_0) \\ f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix}$$

onde

$$J = \frac{1}{2} (E_i)^2 = \frac{1}{2} \left(f(x_i) - \sum_{j=0}^n w_j g_j(x_i) \right)^2$$



e, fazendo

$$\left(\frac{\partial E_i}{\partial w_j} \right)_{j=0, \dots, n} = 0$$

temos,

$$\begin{aligned} \frac{\partial E_i}{\partial w_j} &= \left(f(x_i) - \left(\sum_{j=0}^n w_j g_j(x_i) \right) \right) \left(- \sum_{i=1}^n g_i(x_i) \right) = 0 \\ \Rightarrow - \sum_{i=1}^n f(x_i) g_i(x_i) + \left(\sum_{i=0}^n \sum_{j=0}^n w_j g_j(x_i) g_i(x_i) \right) &= 0 \end{aligned}$$

podemos simplificar a notação da solução por meio do **produto interno**. Ou seja,

$$\Rightarrow \left(\sum_{j=0}^n w_j \langle g_j(\mathbf{x}), g_i(\mathbf{x}) \rangle \right) - \langle f(\mathbf{x}), g_i(\mathbf{x}) \rangle = 0$$

$$\sum_{j=0}^n w_j \langle g_j(\mathbf{x}), g_i(\mathbf{x}) \rangle = \langle f(\mathbf{x}), g_i(\mathbf{x}) \rangle, \quad i = 0, 1, \dots, n$$

de modo que a solução é dada por,

$$\begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{pmatrix} = \begin{pmatrix} \langle g_0(\mathbf{x}), g_0(\mathbf{x}) \rangle & \langle g_0(\mathbf{x}), g_1(\mathbf{x}) \rangle & \cdots & \langle g_0(\mathbf{x}), g_n(\mathbf{x}) \rangle \\ \langle g_1(\mathbf{x}), g_0(\mathbf{x}) \rangle & \langle g_1(\mathbf{x}), g_1(\mathbf{x}) \rangle & \cdots & \langle g_1(\mathbf{x}), g_n(\mathbf{x}) \rangle \\ \vdots & \ddots & \cdots & \vdots \\ \langle g_n(\mathbf{x}), g_0(\mathbf{x}) \rangle & \langle g_n(\mathbf{x}), g_1(\mathbf{x}) \rangle & \cdots & \langle g_n(\mathbf{x}), g_n(\mathbf{x}) \rangle \end{pmatrix}^{-1} \begin{pmatrix} \langle f(\mathbf{x}), g_0(\mathbf{x}) \rangle \\ \langle f(\mathbf{x}), g_1(\mathbf{x}) \rangle \\ \vdots \\ \langle f(\mathbf{x}), g_n(\mathbf{x}) \rangle \end{pmatrix} \quad (6.8)$$

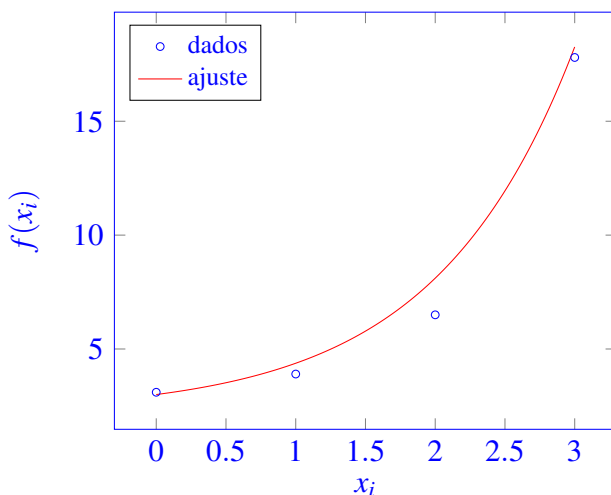


Exemplo 6.6

Encontre uma função na forma $f(x) = w_0 \cos(x) + w_1 e^x$ que melhor se ajuste ao seguinte conjunto de dados discretos:

x_i	0	1	2	3
$f(x_i)$	3,18	3,9	6,5	17,82

Figura 6.7- Gráfico do exemplo 6.4.



$$\langle g_0, g_0 \rangle = \sum_{i=0}^3 \cos^2(x_i) = \cos^2(0) + \cos^2(1) + \cos^2(2) + \cos^2(3) = 2,4452$$

$$\langle g_0, g_1 \rangle = \sum_{i=0}^3 \cos(x_i) e^{x_i} = \cos(0)e^0 + \cos(1)e^1 + \cos(2)e^2 + \cos(3)e^3 = -20,4907$$

pela simetria do produto interno,



$$\langle g_0, g_1 \rangle = \langle g_1, g_0 \rangle = -20,4907$$

$$\langle g_1, g_1 \rangle = \sum_{i=0}^3 e^{2x} = e^0 + e^2 + e^4 + e^9 = 466,4160$$

$$\langle y, g_0 \rangle = \sum_{i=0}^3 y_i \cos(x_i) = y_0 \cos(x_0) + y_1 \cos(x_1) + y_2 \cos(x_2) + y_3 \cos(x_3) = -15,0594$$

$$\langle y, g_1 \rangle = \sum_{i=0}^3 y_i e^{x_i} = y_0 e^{x_0} + y_1 e^{x_1} + y_2 e^{x_2} + y_3 e^{x_3} = 419,7344$$

cuja solução é

$$\begin{pmatrix} w_0 \\ w_1 \end{pmatrix} = \begin{pmatrix} 2,4452 & -20,4907 \\ -20,4907 & 466,4160 \end{pmatrix}^{-1} \begin{pmatrix} -15,0594 \\ 419,7344 \end{pmatrix} = \begin{pmatrix} 2,1880 \\ 0,9960 \end{pmatrix}$$

$$\rightarrow f(x) = 2,1880 \cos(x) + 0,9960 e^x$$

$$RMSE = \sqrt{4 \frac{\|\mathbf{r}\|_2}{4}} = \frac{1}{2} 0,0556 = 0,028.$$

caracterizando o ajuste que pode ser conferido no gráfico da seguinte figura,

FUNÇÕES DE BASE RADIAL - RBF

Seja $\varphi : [0, \infty) \rightarrow \mathbb{R}$, uma função, que quando combinada com uma métrica num espaço vetorial $\|\cdot\| : V \rightarrow \mathbb{R}^+$, de modo que, sobre algum x e $c \in V$, $\varphi(x, c) = \varphi(\|x - c\|)$ é uma funções radial de centro c . Funções radiais, associadas a núcleos radiais são denominadas funções de base radial - RBF se, para um dado conjunto $\{\mathbf{x}_k\}_{k=0}^n$. Os núcleos $\varphi_{x_1}, \varphi_{x_2}, \dots, \varphi_{x_n}$ são linearmente independentes e formam uma base para um espaço de Haar, de modo que a matriz,



$$\begin{pmatrix} \phi(\|x_1 - x_1\|) & \phi(\|x_1 - x_2\|) & \cdots & \phi(\|x_1 - x_n\|) \\ \phi(\|x_1 - x_2\|) & \phi(\|x_2 - x_2\|) & \cdots & \phi(\|x_n - x_2\|) \\ \vdots & \cdots & \ddots & \vdots \\ \phi(\|x_1 - x_n\|) & \phi(\|x_2 - x_n\|) & \cdots & \phi(\|x_n - x_n\|) \end{pmatrix}$$

seja não singular.

A combinação linear de funções de base radial são tipicamente usadas para aproximar funções. Este processo de aproximação também pode ser interpretado como um caso simples de rede neural. Em particular, utilizaremos as funções de base radial do tipo gaussianas. Ou seja,

$$\phi(x, c) = e^{\left(-\frac{1}{2\sigma^2}\|x-c\|^2\right)} \quad (6.9)$$

Assim, o problema Iremos apresentar na presente seção um modelo de ajuste bastante eficaz e estatisticamente robusto. Se o modelo que desejamos ajustar for uma combinação linear de várias funções de base $g_i(x)$ na forma

$$y \approx f(x) = \sum_{i=0}^m w_i \phi_i(x) = w_0 \phi_0(x) + w_1 \phi_1(x) + \cdots + w_m \phi_m(x) \quad (6.10)$$

onde as constantes w 's são desconhecidas, e as ϕ 's são as funções de base em um espaço de funções. Que pode ser escrito como o seguinte sistema de equações lineares,



$$\begin{cases} w_0\phi_0(x_0) + w_1\phi_1(x_0) + w_2\phi_2(x_0) + \cdots + w_m\phi_m(x_0) = f(x_0) \\ w_0\phi_0(x_1) + w_1\phi_1(x_1) + w_2\phi_2(x_1) + \cdots + w_m\phi_m(x_1) = f(x_1) \\ w_0\phi_0(x_2) + w_1\phi_1(x_2) + w_2\phi_2(x_2) + \cdots + w_m\phi_m(x_2) = f(x_2) \\ \vdots \\ w_0\phi_0(x_n) + w_1\phi_1(x_n) + w_2\phi_2(x_n) + \cdots + w_m\phi_m(x_n) = f(x_n). \end{cases}$$

Cuja solução é obtida por meio de,

$$\mathbf{w} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{f}.$$

AJUSTE NÃO LINEAR

Existem algumas funções em que não é possível escrever a priori como uma combinação linear de funções de base, nestes casos se faz necessário algum “artifício matemático” para ser possível ajustar o conjunto de dados com comportamento similar as funções citadas. Como exemplo, temos:

Exemplo 6.7

$$y = \alpha_1 e^{\beta_1 x}$$

$$\ln(y) = \ln(\alpha_1) + \beta_1 x \Rightarrow \begin{cases} f(x) = \ln(y) \\ g_0 = 1 \\ g_1 = x \\ w_0 = \ln(\alpha_1) \rightarrow \alpha_1 = e^{w_0} \\ w_1 = \beta_1 \end{cases}$$

$$y = \alpha_2 x^{\beta_2}$$



$$\log(y) = \log(\alpha_2) + \beta_2 \log(x) \Rightarrow \begin{cases} f(x) = \log(y) \\ g_0 = 1 \\ g_1 = \log(x) \\ w_0 = \log(\alpha_2) \rightarrow \alpha_2 = 10^{w_0} \\ w_1 = \beta_2 \end{cases}$$

$$y = \alpha_3 \frac{x}{\beta_3 + x}$$

$$\frac{1}{y} = \frac{\beta_3}{\alpha_3 x} + \frac{1}{\alpha_3} \Rightarrow \begin{cases} f(x) = \frac{1}{y} \\ g_0 = \frac{1}{x} \\ g_1 = 1 \\ w_0 = \frac{\beta_3}{\alpha_3} \\ w_1 = \frac{1}{\alpha_3} \end{cases}$$

Podemos de certa forma “linearizar” estas funções de forma a usar a regressão linear por mínimos quadrados definida na equação (6.2).

Exemplo 6.8

Um experimento no laboratório de física gerou os pontos dados na seguinte tabela. Obtenha uma ajuste do tipo $y = \alpha e^{\beta x}$ para o conjunto.

x_i	-2	-1	0	1	2	3
y_i	0,05	0,15	0,4	1,1	2,3	7,1



$$\left\{ \begin{array}{l} \langle g_0, g_0 \rangle = \sum_{i=0}^5 1 = 6 \\ \langle g_0, g_1 \rangle = \langle g_1, g_0 \rangle = \sum_{i=0}^5 x_i = 3 \\ \langle g_1, g_1 \rangle = \sum_{i=0}^5 x_i^2 = 19 \\ \langle f, g_0 \rangle = \sum_{i=0}^5 \ln(y_i) = -2,9208 \quad \langle f, g_1 \rangle = \sum_{i=0}^5 \ln(y_i)x_i = 15,5299 \end{array} \right.$$

cuja solução é obtida por meio de

$$\begin{pmatrix} w_0 \\ w_1 \end{pmatrix} = \begin{pmatrix} 6 & 3 \\ 3 & 19 \end{pmatrix}^{-1} \begin{pmatrix} -2,9208 \\ 15,5299 \end{pmatrix} = \begin{pmatrix} -0,9722 \\ 0,9708 \end{pmatrix} \rightarrow \alpha = e^{w_0} = 0,3782$$

com

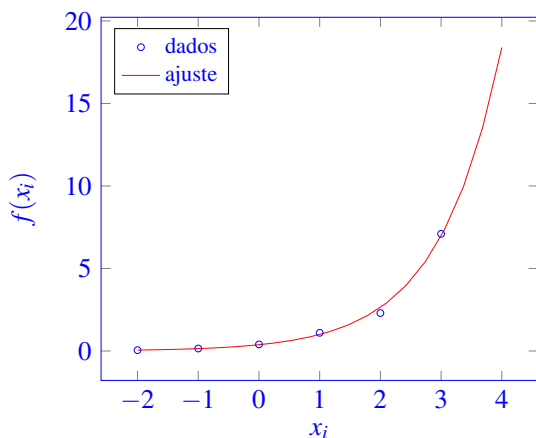
$$f(x) = \alpha e^{\beta x} = 0,3782 e^{0,9708x}$$

e

$$RMSE = \frac{\sqrt{6}}{6} 0,3786 = 0,1545$$



Figura 6.8 – Gráfico do exemplo 6.10.



AJUSTE TRIGONOMÉTRICO

Agora iremos abordar um tipo especial de aproximação ou ajuste utilizando mínimos quadrados denominada aproximação de Fourier. Para isso iremos considerar a equação (6.10), escrita da seguinte forma,

$$f(x) = \frac{a_0}{2} + a_1 \cos(x) + a_2 \cos(2x) + \dots + a_n \cos(nx) + b_1 \sin(x) + b_2 \sin(2x) + \dots + b_m \sin(mx). \quad (6.11)$$

Num subespaço do $C([0, 2\pi])$ com funções de base,

$$\wp = \{1, \cos(x), \cos(2x), \dots, \cos(nx), \sin(x), \sin(2x), \dots, \sin(mx)\}$$

que são ortogonais, tal como mostrado no exercício A.3 (anexo A). No entanto, no mesmo exercício, vimos que normalizando cada função da base \wp , podemos obter uma base ortonormal,



$$B = \{\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_n, \mathbf{v}_{n+1}, \dots, \mathbf{v}_{2n}\}$$

$$= \left\{ \frac{1}{\sqrt{2\pi}}, \frac{1}{\sqrt{\pi}}\cos(x), \dots, \frac{1}{\sqrt{\pi}}\cos(nx), \frac{1}{\sqrt{\pi}}\sen(x), \dots, \frac{1}{\sqrt{\pi}}\sen(mx) \right\}.$$

Pela teoria de aproximação de funções (proposição A.1), podemos escrever cada função nesse subespaço como,

$$f(x) = \langle f, \mathbf{v}_0 \rangle \mathbf{v}_0 + \langle f, \mathbf{v}_1 \rangle \mathbf{v}_1 + \dots + \langle f, \mathbf{v}_n \rangle \mathbf{v}_n. \quad (6.12)$$

Para que a equação (6.11) seja igual a (6.12). É necessário que,

$$\frac{a_0}{2} = \langle f, \mathbf{v}_0 \rangle \mathbf{v}_0 \rightarrow a_0 = \frac{2}{\sqrt{2\pi}} \int_0^{2\pi} f(x) \frac{1}{\sqrt{2\pi}} dx \rightarrow a_0 = \frac{1}{\pi} \int_0^{2\pi} f(x) dx$$

$$a_1 \cos(x) = \langle f, \mathbf{v}_1 \rangle \mathbf{v}_1 \rightarrow a_1 \cos(x) = \frac{\cos(x)}{\sqrt{\pi}} \int_0^{2\pi} f(x) \frac{\cos(x)}{\sqrt{\pi}} dx \rightarrow a_1 = \frac{1}{\pi} \int_0^{2\pi} f(x) \cos(x) dx$$

$$\vdots$$

$$a_n \cos(nx) = \langle f, \mathbf{v}_n \rangle \mathbf{v}_n \rightarrow a_n = \frac{1}{\pi} \int_0^{2\pi} f(x) \cos(nx) dx$$

$$b_1 \sen(x) = \langle f, \mathbf{v}_{n+1} \rangle \mathbf{v}_{n+1} \rightarrow b_1 = \frac{1}{\pi} \int_0^{2\pi} f(x) \sen(x) dx$$

$$\vdots$$

$$b_m \sen(mx) = \langle f, \mathbf{v}_{2n} \rangle \mathbf{v}_{2n} \rightarrow b_m = \frac{1}{\pi} \int_0^{2\pi} f(x) \sen(mx) dx.$$

A função $f(x)$ é a aproximação trigonométrica de Fourier no intervalo $[0 ; 2\pi]$. Com coeficientes

$$\left\{ \begin{array}{l} a_0 = \frac{1}{\pi} \int_0^{2\pi} f(x) dx \\ a_k = \frac{1}{\pi} \int_0^{2\pi} f(x) \cos(kx) dx, \quad k = 1, 2, \dots, n \\ b_k = \frac{1}{\pi} \int_0^{2\pi} f(x) \sen(kx) dx, \quad k = 1, 2, \dots, m \end{array} \right. \quad (6.13)$$



Exemplo 6.9

Encontre a aproximação de Fourier de n-ésima ordem da função $f(x) = x$ no intervalo $[0, 2\pi]$.

$$f(x) = \frac{a_0}{2} + \sum_{k=1}^n (a_k \cos(kx) + b_k \sin(kx))$$

$$a_0 = \frac{1}{\pi} \int_0^{2\pi} x dx = \frac{x^2}{2} \Big|_0^{2\pi} = \frac{4\pi^2}{2\pi} = 2\pi$$

$$a_k = \frac{1}{\pi} \int_0^{2\pi} x \cos(kx) dx = \frac{1}{k\pi} \left[\underbrace{(x \sin(kx)) \Big|_0^{2\pi}}_0 - \underbrace{\int_0^{2\pi} \sin(kx) dx}_0 \right] = 0$$

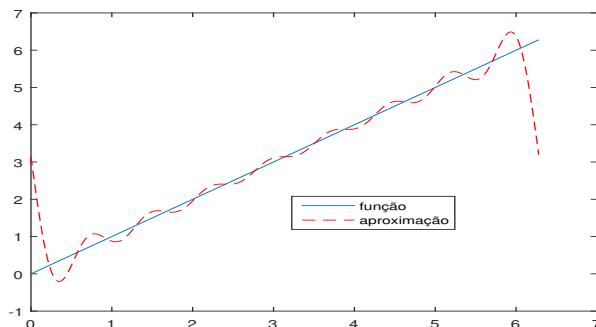
$$b_k = \frac{1}{\pi} \int_0^{2\pi} x \sin(kx) dx = -\frac{1}{k\pi} \left[\underbrace{(x \cos(kx)) \Big|_0^{2\pi}}_{2\pi} - \underbrace{\int_0^{2\pi} \cos(kx) dx}_0 \right] = -\frac{2}{k}$$

cuja aproximação de ordem n-ésima é dada por,

$$\begin{aligned} &= \pi - 2 \sin(x) + \frac{2}{2} \sin(2x) + \frac{2}{3} \sin(3x) + \frac{2}{4} \sin(4x) + \dots + \frac{2}{m} \sin(mx) \\ &= \pi - 2 \left(\sin(x) + \frac{1}{2} \sin(2x) + \frac{1}{3} \sin(3x) + \frac{1}{4} \sin(4x) + \dots + \frac{1}{m} \sin(mx) \right) \\ f(x) &\approx \pi - 2 \sum_{k=1}^m \frac{\sin(kx)}{k} \end{aligned}$$



Figura 6.9 – aproximação da função $f(x) = x$, com $m = 8$



R função par $f(x) = f(-x)$

função ímpar $f(-x) = -f(x)$

Se a função dada for uma função par e como $\text{sen}(x)$ é uma função ímpar, temos:

$$\int_0^{2\pi} f(x) \text{sen}(x) dx = 0, \quad (6.14)$$

portanto,

$$f(x) \approx \frac{a_0}{2} + \sum_{k=1}^n a_k \cos(kx). \quad (6.15)$$

Caso contrário, ou seja, $f(x)$ for ímpar, temos:

$$f(x) \approx \sum_{k=1}^m b_k \text{sen}(kx). \quad (6.16)$$



Exemplo 6.10

Obter uma aproximação de terceira ordem para a função:

$$f(x) = |x|, \quad -\pi \leq x < \pi.$$

Como a função $f(x) = |x|$ é uma função par, a aproximação será dada pela expressão (6.15), com

$$\begin{aligned} a_0 &= \frac{2}{\pi} \int_0^{\pi} x dx = \frac{2}{\pi} \left[\frac{x^2}{2} \right]_0^{\pi} = \pi \\ a_k &= \frac{2}{\pi} \int_0^{\pi} x \cos(kx) dx = \frac{2}{\pi} \left[\underbrace{\left[-\frac{x \operatorname{sen}(kx)}{k} \right]_0^{\pi}}_0 + \frac{1}{k} \int_0^{\pi} \operatorname{sen}(kx) dx \right] \\ &= \frac{2}{k^2 \pi} [\cos(kx)]_0^{\pi} = -\frac{4}{k^2 \pi} \\ f(x) &\approx \frac{\pi}{2} - \frac{4}{\pi} \sum_{k=1}^n \frac{\cos((2k-1)x)}{(2k-1)^2} \end{aligned}$$

Uma outra forma, geralmente a mais conveniente, pode ser obtida por meio das fórmulas de Euler,

$$\begin{cases} \cos(kx) = \frac{(e^{ikx} + e^{-ikx})}{2} \\ e \\ \operatorname{sen}(kx) = \frac{(e^{ikx} - e^{-ikx})}{2i} \end{cases} \quad (6.17)$$

onde $i^2 = -1$ é a unidade imaginária. Substituindo as equações (6.17) na expressão (6.11). Obtemos (considerando, $m = n$),



$$f(x) = \frac{a_0}{2} + \sum_{k=1}^n \left(a_k \frac{(e^{ikx} + e^{-ikx})}{2} + b_k \frac{(e^{ikx} - e^{-ikx})}{2i} \right)$$

$$f(x) = \frac{1}{2} \left[a_0 + \sum_{k=1}^n \left[a_k (e^{ikx} + e^{-ikx}) - ib_k (e^{ikx} - e^{-ikx}) \right] \right]$$

$$f(x) = \frac{1}{2} \left[a_0 + \sum_{k=1}^n \left[(a_k - ib_k) e^{ikx} + (a_k + ib_k) e^{-ikx} \right] \right]$$

fazendo,

$$c_0 = \frac{a_0}{2}, c_k = \frac{a_k - ib_k}{2} \text{ e } c_k = \frac{a_k + ib_k}{2},$$

podemos reescrever a expressão anterior como,

$$f(x) = \sum_{k=-n}^n c_k e^{ikx} \quad (6.18)$$

que corresponde uma outra forma de aproximação trigonométrica utilizando variáveis complexas graças a presença da exponencial complexa. O foco do nosso estudo têm sido funções periódicas com período fundamental 2π , que são totalmente definidas por seus valores num intervalo de $[-\pi, \pi]$, ou $[0, 2\pi]$. Se uma função de x possui período Δx , então podemos dividir o intervalo considerando, por exemplo, $t = 2\pi \frac{n}{N} x$, $n = 0, 1, \dots, N - 1$ que converte a função em uma função na variável t com período 2π . A função pode ter valores complexos, pois a função exponencial complexa é conveniente para manipulações. Nosso problema se resume então a encontrar $c_k \in \mathbb{C}$ tal que, se a série na equação (6.18) converge para



$f(x)$ no intervalo $0 \leq x \leq 2\pi$. Para algum inteiro m , multiplicando (6.18) por e^{-imx} , obtemos:

$$f(x)e^{-imx} = \sum_{k=-n}^n c_k e^{i(k-m)x} \rightarrow \int_0^{2\pi} f(x)e^{-imx} dx = \sum_{k=-n}^n \int_0^{2\pi} c_k e^{i(k-m)x} dx;$$

no entanto,

$$\int_0^{2\pi} e^{i(k-m)x} dx = \begin{cases} \left. \frac{e^{i(n-m)x}}{i(n-m)} \right|_0^{2\pi} = 0 & \text{se } k \neq m \\ \left. x \right|_0^{2\pi} = 2\pi & \text{se } k = m \end{cases}$$

onde,

$$c_m = \frac{1}{2\pi} \int_0^{2\pi} f(x)e^{-imx} dx, \quad (6.19)$$

que é válida se a série descrita em (6.18) for convergente.

Definição 6.1

Se f possui período 2π e é integrável no intervalo de $[-\pi; \pi]$, a série

$$\sum_{n=-\infty}^{\infty} c_n e^{int}$$

com coeficientes obtidos via (6.4) é denominada a série de Fourier da f ; e os números complexos c_n são os coeficientes de Fourier da função.

Geralmente, se a função $f(t)$ possui período T , a série de Fourier pode ser escrita como,

$$\sum_{n=-\infty}^{\infty} c_n e^{2\pi i \frac{nt}{T}},$$



e os coeficientes de Fourier,

$$c_n = \frac{1}{T} \int_{-\frac{T}{2}}^{\frac{T}{2}} f(t) e^{-2\pi i \frac{nt}{T}} dt.$$

No entanto, a maioria das funções de uso na engenharia, são representadas por meio de um conjunto finito de valores discretos. Quando não, os dados são coletados ou convertidos do contínuo para o discreto por meio de conversores. Neste caso, representamos a função discreta como,

$$F_k(n) = \sum_{n=0}^{N-1} f(n) e^{-2\pi i \frac{nk}{N}}, \quad k = 0, 1, \dots, N-1.$$

cuja inversa é dada por

$$f(n) = \frac{1}{N} \sum_{k=0}^{N-1} F_k(n) e^{2\pi i \frac{nk}{N}}, \quad n = 0, 1, \dots, N-1.$$

O subscrito n é usado para designar os instantes discretos nos quais as medidas foram tomadas. Logo, $f(n)$ designa o valor da função definida no contínuo $f(t)$ no instante t_n .

EXERCÍCIOS PROPOSTOS

Exercício 6.1

Considere o seguinte funcional de custo

$$J = \frac{1}{2} \sum_{i=1}^n (y_i - f(x_i))^2$$



onde $f(x) = ax + b\sin(x)$ e $(x_i; y_i)$ é um conjunto de pontos no plano. Para minimizar o funcional de custo, tome as derivadas parciais em relação as parâmetros a e b , monte o SEL e encontre sua solução para os pares ordenados: (1; 68), (1, 2; 2, 05) e (2; 2, 51).

Exercício 6.2

Suponha que durante um certo experimento num laboratório de Física, foi obtido o seguinte conjunto de dados,

t	0	1	2	3	4	5
$x(t)$	-0,3698	2,5298	6,9162	8,4871	13,9238	14,9146

O pesquisador acredita numa relação linear entre os dados, mas ao fazer uma análise prévia (por meio do gráfico), não é possível afirmar a existência de uma linha perfeita devido ao erro de medição. Portanto, se faz necessário uma ajuste nos dados por meio de algum método. Para isso, iremos considerar o modelo de ajuste $x(t) = at + b$.

Exercício 6.3

É conhecido que a queda de voltagem através de um indutor segue a lei de Faraday,

$$V_L = L \cdot \frac{di}{dt}$$



onde V_L é a queda de tensão no indutor (em volts), L é a indutância (em henrys) e i a corrente (em ampères). Use os seguintes dados para fazer uma estimativa do valor da indutância,

$\frac{di}{dt}, \frac{A}{s}$	1	2	4	6	8	10
V_L, V	5,5	12,5	17,5	32	38	49

Exercício 6.4

Aproxime a seguinte função no intervalo de $[-\pi; \pi]$.

$$h(x) = \sum_{n=0}^{\infty} (-1)^n \frac{x^n}{n!} + \sum_{n=0}^{\infty} \frac{x^n}{n!}$$

no intervalo de $[-\pi; \pi]$.

Exercício 6.5

Num experimento em laboratório, foi analisado o efeito da temperatura na resistência de um metal. Os valores obtidos estão representados na seguinte tabela:

$T(^{\circ}C)$	20,5	32,7	51,0	73,2	95,7
$R(\Omega)$	765	826	873	942	1032



encontre uma função que melhor se ajuste ao conjunto de pontos e estime a resistência para uma temperatura de 100°C .

Exercício 6.6

O pH em um reator varia senoidalmente no decorrer de um dia. Utilize regressão por mínimos quadrados para ajustar uma curva senoidal ao conjunto de dados,

$t(h)$	0	2	4	5	7	9	12	15	20	22	24
pH	7,6	7,2	7	6,5	7,5	7,2	8,9	9,1	8,9	7,9	7

Use o ajuste obtido para determinar a média, a amplitude e o instante de pH máximo. Observe que o período é de $24h$.

Exercício 6.7

Após serem efetuadas medições num gerador elétrico, foram obtidos os seguintes valores,

$I(A)$	1,58	2,15	4,8	4,9	3,12	3,01
$V(v)$	210	180	150	120	60	30

- ajuste os dados por uma função polinomial de grau adequado;



- estime o valor da tensão obtida no voltímetro para um valor de 3,05A.

Exercício 6.8

Um experimento com um circuito RC é usado para determinar a capacitância de um dado capacitor. No circuito, um resistor é ligado em série com um capacitor e uma bateria, de modo que os dados medidos no experimento são dados pela seguinte tabela,

$t(s)$	2	4	6	8	10	12
$V_r(V)$	4,41	1,62	0,6	0,22	0,08	0,03

teoricamente, a tensão (em função do tempo) no resistor é dada por,

$$V_r = V e^{-t/RC}$$

- encontrar a função que melhor ajusta os dados;
- determinar a capacitância (do capacitor), considerando $R = 5\Omega$.

Exercício 6.9

Os seguintes dados mostram a relação entre a viscosidade de um determinado óleo e a temperatura. Use regressão não linear para ajustar uma equação de potência para esses dados,



T	26,67	93,33	148,89	315,56
μ	1,35	0,085	0,012	0,00075

onde μ representa a viscosidade e T a temperatura.

Exercício 6.10

A função

$$y = e^{(a-b.e^{-ct})} = \exp(a - b.\exp(-c.t)), \quad b, c > 0; a \in \mathbb{R}$$

aparece frequentemente em inúmeras áreas da pesquisa. Em particular, em ciências atuariais para especificar a taxa de mortalidade, na medicina para modelar crescimento em tumores, na biologia em modelagem de crescimento em culturas e sistemas, na ecologia, no marketing, etc.

Observando o crescimento de um certo tipo de animal ao longo de 9 semanas, foi obtida a seguinte tabela, relacionando o tempo ao peso médio(em kg):

t_k	1	2	3	4	5	6	7	8	9
y_k	0,320	0,3238	0,5126	0,4602	0,7897	0,6906	0,9981	1,07	1,10

Obtenha os valores de a, b que melhor aproxime o conjunto de pontos pela função dada, considerando $c = 0,2385$.



Exercício 6.11

Seja o seguinte conjunto de pontos,

x_k	-1,0	-0,6	-0,2	0,2	0,6	1,0
y_k	-0,9602	0,6405	-0,2013	-0,1647	0,3449	-0,3201
c_i	-1,5	-0,5	0,0	0,4	1,0	

encontre a função

$$f(x) = a_0\phi_0(x) + a_1\phi_1(x) + a_2\phi_2(x) + a_3\phi_3(x) + a_4\phi_4(x)$$

com $\phi_i(x) = \exp(-(2\sigma^2)^{-1}(xk - c_i)^2)$, $i = 0, 1, \dots, 4$ (e , σ definido pelo usuário) que melhor ajusta o conjunto de pontos. De modo que $RMSE \leq 0,03$.

Exercício 6.12

Os resultados de um experimento em um túnel de vento sobre o fluxo de ar na ponta da asa de um avião forneceu os seguintes dados:

R/C	1	2	3	4	5	6	7	8	9
V_θ/V_∞	0,147	0,357	0,641	0,980	1,358	1,758	2,159	2,549	3,251

onde R é a distância do núcleo do vórtice, C é a corda de perfil de asa da aeronave, V_θ é a velocidade tangencial



do vórtice e V_∞ é a velocidade de fluxo livre da aeronave. Seja $x = R/C$ e $y = V_0/V_\infty$. Obtenha uma função que melhor se ajusta ao conjunto de pontos ou ajuste-os por meio da seguinte função.

$$f(x) = \frac{\alpha}{x} \left(1 - e^{-\beta x^2}\right).$$

Exercício 6.13

Obter uma aproximação de quarta ordem para a função:

$$f(x) = \begin{cases} -1, & -\pi \leq x < 0 \\ 1, & 0 \leq x < \pi \end{cases}$$

Exercício 6.14

Suponha a medição da temperatura ao longo de 19 horas, representadas pelo seguinte conjunto de dados,



t (horas)	$T(t)$ – temperatura ($^{\circ}\text{C}$)
0	19.4
1	18.9
2	18.2
3	17.9
4	17.1
5	15.8
6	12.5
7	14.6
8	15.1
9	14.6
10	14.3
11	16.7
12	18
13	19.7
14	21.6
15	23
16	24
17	23.3
18	23.8
19	23.6

Obtenha um modelo da forma $T \approx f(t) = a_0g_0(t) + a_1g_1(t) + a_2g_2(t) + a_3g_3(t)$ que melhor ajuste a variação da temperatura ao longo do dia. Onde $g_i(t)$ com $i = 0, \dots, 3$ são funções de base radial.



Exercício 6.15

Faça uma aproximação trigonométrica das seguintes funções:

$$f(x) = \begin{cases} 0, & -\pi \leq x < 0 \\ x, & 0 \leq x < \pi \end{cases} \quad n = 6$$

$$f(x) = \begin{cases} -1, & -\pi \leq x < 0 \\ 3, & 0 \leq x < \pi \end{cases} \quad n = 6$$

$$f(x) = \frac{x}{2}, x \in [-\pi, \pi[, \quad n = 5$$



INTERPOLAÇÃO

Interpolação polinomial geralmente é o início de algum processo numérico. Geralmente, são utilizados como blocos de construção para algoritmos mais complexos em diferenciação, integração, solução de equações diferenciais, teoria da aproximação, bem como outras aplicações. A interpolação é um caso especial de ajuste de curvas ou aproximação. Por exemplo, assim como ajuste de curvas, dado um conjunto de pontos,

$$(x_i; y_i)_{i=0}^n$$

interpolar equivale a encontrar uma função $f(x)$ que passe exatamente pelos pontos dados, satisfazendo

$$f(x_i) = y_i, \quad i = 0, 1, \dots, n.$$

Geralmente, representamos a função interpolante, como uma combinação linear de funções de base, linearmente independentes.

Ou seja,

$$f(x) = \sum_{j=0}^n w_j \phi_j(x) = w_0 \phi_0(x) + w_1 \phi_1(x) + \dots + w_n \phi_n(x),$$

assim, como num problema de ajuste de curvas,

$$\{w_j\}_{j=0}^n$$

são os parâmetros desconhecidos do modelo e

$$\{\phi_j(x)\}_{j=0}^n$$



funções de base pré-definidas. Neste curso, iremos considerar a interpolação polinomial.

Teorema 7.1

Existência do polinômio interpolador. Sejam x_0, x_1, \dots, x_n pontos distintos, então para os valores y_0, y_1, \dots, y_n existe um único polinômio p de grau menor ou igual a n tal que $p(x_i) = y_i$ para $0 \leq i \leq n$.

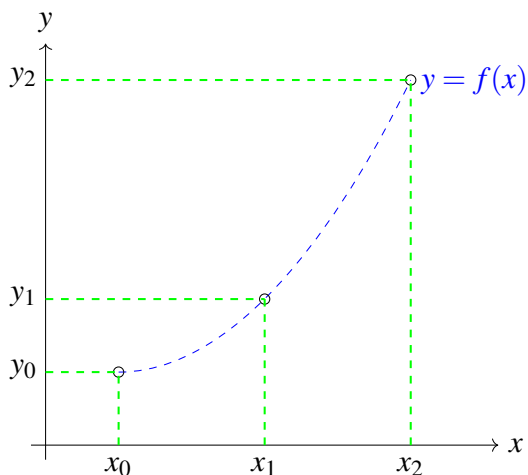
POLINÔMIOS INTERPOLADORES DE LAGRANGE

Seja a função $y = f(x)$ definida num intervalo $[x_0, x_2]$, cujo gráfico é dado pela figura 7.1. Aparentemente o gráfico consiste de uma função quadrática. Portanto, devemos aproximar nossa função por meio de um simples polinômio de ordem 2. Ou seja,

$$p_2(x) = w_0 + w_1x + w_2x^2.$$



Figura 7.1 – aproximação da função quadrática



O método de interpolação consiste em no intervalo dado $[x_0, x_2]$, escolher um ponto interior x_1 , cujas imagens são:

$$\begin{cases} y_0 = f(x_0) \\ y_1 = f(x_1) \\ y_2 = f(x_2). \end{cases}$$

Iremos construir o polinômio (7.1) tal que os pontos x_0, x_1, x_2 possuam imagem na função quadrática original. Ou seja,

$$\begin{cases} p_2(x_0) = y_0 \\ p_2(x_1) = y_1 \\ p_2(x_2) = y_2. \end{cases}$$

Para auxiliar na resolução do problema, vamos começar construindo um polinômio auxiliar $L_i(x_j)$ de segunda ordem, satisfazendo as seguintes condições,



$$L_i(x_j) = \begin{cases} 1 & \text{se } i = j \\ 0 & \text{se } i \neq j. \end{cases}$$

Por exemplo, $L_0(x_0) = 1$, $L_0(x_1) = 0$ e $L_0(x_2) = 0$. Portanto,

$$L_0(x) = A(x - x_1)(x - x_2).$$

Mas,

$$L_0(x_0) = A(x_0 - x_1)(x_0 - x_2) = 1 \rightarrow A = \frac{1}{(x_0 - x_1)(x_0 - x_2)}.$$

Substituindo na equação (7.2), obtendo:

$$L_0(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)}.$$

De maneira similar podemos obter,

$$L_1(x) = \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)}$$

$$L_2(x) = \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)}.$$

Portanto o polinômio desejado é dado pela seguinte expressão:

$$\begin{aligned} p_2(x) &= y_0 L_0(x) + y_1 L_1(x) + y_2 L_2(x) \\ &= y_0 \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} + y_1 \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} + y_2 \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)}. \end{aligned}$$



Exemplo 7.1

Interpolar a curva $y = \text{sen}(x)$ no intervalo $[0, \pi]$, por um polinômio de grau 2.

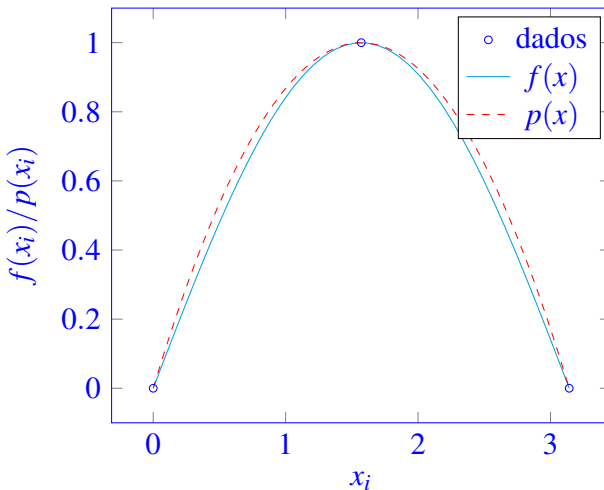
Escolhendo $x_1 = \pi/2$, temos que,

$$\begin{cases} y_0 = \text{sen}(0) = 0 \\ y_1 = \text{sen}(\frac{\pi}{2}) = 1 \\ y_2 = \text{sen}(\pi) = 0 \end{cases}$$

que resume nossa solução ao seguinte polinômio,

$$p_2(x) = y_1 L_1(x) = \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} = \frac{(x-0)(x-\pi)}{(\frac{\pi}{2}-0)(\frac{\pi}{2}-\pi)} = \frac{4}{\pi^2}x(\pi-x).$$

Figura 7.2 – Interpolação da função do exemplo 7.1.



Generalizando, vamos denotar um interpolador polinomial de grau n , por



$$f(x) \approx p_n(x) = \sum_{j=0}^n w_j x^j = w_0 + w_1 x + w_2 x^2 + \dots + w_n x^n.$$

para $n + 1$ pontos de dados. Na interpolação devemos calcular os coeficientes desconhecidos w_0, w_1, \dots, w_n de modo que

$$p_n(x_i) = y_i, \quad i = 0, 1, \dots, n.$$

Assumiremos, inicialmente que $x_i \neq x_j$ sempre que $i \neq j$ e sejam, os $n + 1$ polinômios $p_i(x)$ de grau n escritos na sua forma fatorada, através do seguinte SEL:

$$\left\{ \begin{array}{l} p_0(x) = (x - x_1)(x - x_2) \dots (x - x_n) \\ p_1(x) = (x - x_0)(x - x_2) \dots (x - x_n) \\ p_2(x) = (x - x_0)(x - x_1) \dots (x - x_n) \\ \vdots \\ p_n(x) = (x - x_0)(x - x_1) \dots (x - x_{n-1}) \end{array} \right.$$

Ou de forma simplificada

$$p_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n (x - x_j), \quad i = 0, 1, \dots, n$$

que possuem as seguintes propriedades:

$$\left\{ \begin{array}{l} p_i(x_i) \neq 0 \quad \forall i \\ p_i(x_j) = 0 \quad \forall j \neq i. \end{array} \right.$$

e, são conhecidos como **polinômios de Lagrange**. Como o polinômio $p(x)$ que desejamos encontrar é de grau n e contém os pontos $\{(x_i, y_i)\}_{i=0}^n$ podemos escrevê-lo como uma combinação linear dos polinômios $p_i(x)$, $i = 0, 1, \dots, n$. Então,



$$p_n(x) = w_0 + w_1 p_1(x) + \cdots + w_n p_n(x)$$

ou

$$p_n(x) = \sum_{i=0}^n w_i p_i(x)$$

seja

$$p_n(x_m) = \sum_{i=0}^n w_i p_i(x_m) = w_0 p_0(x_m) + w_1 p_1(x_m) + \cdots + w_m p_m(x_m) + \cdots + w_n p_n(x_m)$$

pela propriedade dos polinômios de Lagrange, temos que $p_i(x_m) = 0$ e $p_m(x_m) \neq 0$ de modo que,

$$p_n(x_m) = w_m p_m(x_m)$$

daí temos

$$w_m = \frac{p_n(x_m)}{p_m(x_m)}$$

como consideramos inicialmente que $p_n(x_i) = y_i$, vale a relação:

$$w_i = \frac{y_i}{p_i(x_i)}$$

Concluindo, temos:

$$p_n(x) = \sum_{i=0}^n \frac{y_i}{p_i(x_i)} p_i(x)$$

ou

$$p_n(x) = \sum_{i=0}^n y_i \frac{p_i(x)}{p_i(x_i)}$$



com isso teremos o polinômio interpolador de Lagrange:

$$p_n(x) = \sum_{i=0}^n y_i \prod_{\substack{j=0 \\ j \neq i}}^n \frac{(x - x_j)}{(x_i - x_j)}$$

Exemplo 7.2

Determinar o polinômio interpolador de Lagrange para a função conhecida pelos seguintes pontos e determinar o valor $p(0, 3)$:

i	0	1	2	3
x_i	0	0,2	0,4	0,5
y_i	0	2,0	4,06	5,12

$$L_1(x) = \frac{x(x-0,4)(x-0,5)}{0,2(0,2-0,4)(0,2-0,5)} = \frac{250}{3}x(x-0,4)(x-0,5)$$

$$L_2(x) = \frac{x(x-0,2)(x-0,5)}{0,4(0,4-0,2)(0,4-0,5)} = -125x(x-0,2)(x-0,5)$$

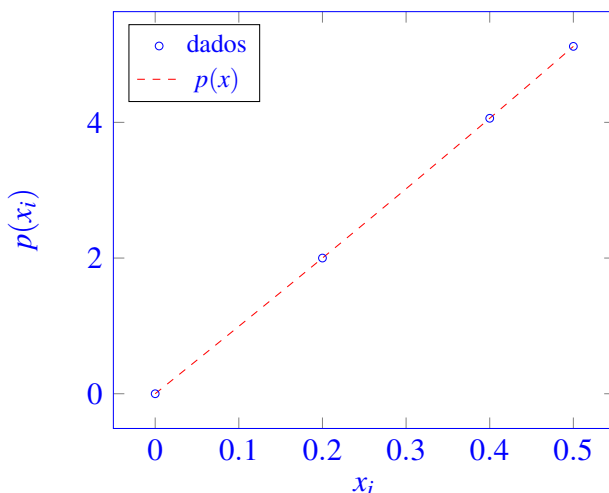
$$L_3(x) = \frac{x(x-0,2)(x-0,4)}{0,5(0,5-0,2)(0,5-0,4)} = \frac{200}{3}x(x-0,2)(x-0,4)$$

de modo que o polinômio interpolador é dado por:

$$p_3(x) = \frac{500}{3}x(x-0,4)(x-0,5) - \frac{1015}{2}x(x-0,2)(x-0,5) + \frac{1024}{3}x(x-0,2)(x-0,4)$$



Figura 7.3 – Interpolação da função do exemplo 7.2.



Teorema 7.2

Erro de truncamento $E(x)$. Seja $f(x)$ uma função definida e $(n + 1)$ vezes diferenciável num intervalo $[a, b]$. Sejam x_0, x_1, \dots, x_n , $(n + 1)$ pontos distintos no intervalo. Se $p(x)$ interpola $f(x)$ nestes pontos, então o erro cometido $E(x)$ é dado por,

$$E(x) = f(x) - p(x) \leq \frac{|\psi(x)|}{(n + 1)!} M$$

em que $M = \max[f^{(n+1)}(x)]$, $x \in [a, b]$, e



$$\psi(x) = \prod_{i=0}^n (x - x_i)$$

Teorema 7.3

Erro de truncamento $E(x)$ para pontos equidistantes.

Seja $f(x)$ uma função definida e $(n + 1)$ vezes diferenciável num intervalo $[a; b]$. Sejam x_0, x_1, \dots, x_n , $(n + 1)$ pontos distintos equidistantes no intervalo. Se $p(x)$ interpola $f(x)$ nestes pontos, então o erro cometido $E(x)$ é dado por,

$$E(x) = f(x) - p(x) \leq \left| \frac{h^2}{8} M \right|$$

em que $M = \max[f^{(2)}(x)], x \in [a, b]$

Exemplo 7.3

Considere a função $f(x) = \cos(x)$, tabelada nos pontos,

x_i	0,2	0,4	0,6
y_i	0,9801	0,9211	0,8253

- Determine o polinômio interpolador de Lagrange, calcule o valor da função no ponto $x = 0,5$ e o limite superior do erro.
- Qual deve ser a amplitude do intervalo a ser considerado na função $f(x) = 1/(1+x)$ no intervalo $[0, 2]$, de modo que utilizando um polinômio interpolador de primeira ordem, apresente um erro menor ou igual a 0,0001?



A desvantagem da abordagem por meio do polinômio interpolador de Lagrange é que cada polinômio $L_i(x)$ depende de todo o conjunto de dados. Portanto, quando se necessita substituir, acrescentar ou remover dados (nem que seja um único), devemos recalculá-los todos os polinômios $L_i(x)$. Uma forma utilizada para contornar tal problema é o que iremos abordar na próxima seção.

POLINÔMIOS INTERPOLADORES POR DIFERENÇAS DIVIDIDAS DE NEWTON

Considere o seguinte polinômio interpolador de grau n

$$p_n(x) = a_0 + a_1(x-x_0) + a_2(x-x_0)(x-x_1) + \cdots + a_n(x-x_0)(x-x_1)\cdots(x-x_{n-1}) = \sum_{j=0}^n a_j g_j(x) \quad (7.3)$$

onde, $g_0(x) = 1$ e

$$g_j(x) = (x-x_0)(x-x_1)\cdots(x-x_{n-1}) = \prod_{s=0}^{n-1} (x-x_s), \quad j = 1, 2, \dots, n-1.$$

Os coeficientes a_j são calculados de maneira recorrente por meio de diferenças divididas.

Definição 7.1

Seja $f(x)$ uma função contínua, $(n+1)$ vezes diferenciável e definida em x_0, x_1, \dots, x_n ($n+1$) pontos distintos num intervalo $[a, b]$.

Definimos **diferença dividida de ordem zero** de uma função $f(x)$ definida nos pontos $x_i, i = 0, 1, \dots, n$ por:

$$f[x_i] = f(x_i), \quad i = 0, 1, \dots, n$$

e, ao aproximarmos a derivada primeira pela reta secante, temos a **diferença dividida de ordem um** (ou primeira ordem), ou seja,



$$f[x_i, x_{i+1}] = \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i}, \quad i = 0, 1, \dots, n-1$$

e, finalmente, temos a **diferença dividida de ordem n** de uma função $f(x)$ definida nos pontos $x_i, i = 0, 1, \dots, n$ por:

$$f[x_0, x_1, \dots, x_n] = \frac{f[x_1, x_2, \dots, x_n] - f[x_0, x_1, \dots, x_{n-1}]}{x_n - x_0}$$

Considerando $x = x_i$ e usando a condição de que $p_n(x_i) = f(x_i), i = 0, 1, \dots, n-1$, no polinômio definido em (7.3) temos que,

$$p_n(x_0) = a_0 + 0 + \dots + 0 = f(x_0) \longrightarrow a_0 = f(x_0)$$

que corresponde a diferença dividida de ordem zero ($f[x_0]$). Para $x = x_1$, temos

$$p_n(x_1) = a_0 + a_1(x_1 - x_0) + 0 + \dots + 0 = f(x_1) \longrightarrow a_1 = \frac{f(x_1) - f(x_0)}{x_1 - x_0},$$

correspondendo a diferença dividida de ordem um $f[x_0, x_1]$. Em seguida, para $x = x_2$

$$p_n(x_2) = a_0 + a_1(x_2 - x_0) + a_2(x_2 - x_0)(x_2 - x_1) + 0 + \dots + 0 = f(x_2) \longrightarrow a_2 = \frac{f(x_2) - a_0 - a_1(x_2 - x_0)}{(x_2 - x_0)(x_2 - x_1)}$$

$$= \frac{\frac{f(x_2) - f(x_1)}{x_2 - x_1} - \frac{f(x_1) - f(x_0)}{x_1 - x_0}}{x_2 - x_0} = \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0},$$

e, assim, sucessivamente obtemos uma expressão para a_n que corresponderá a diferença finita de ordem n , representada por meio do seguinte teorema,



Teorema 7.4

Diferenças divididas. Seja $f(x)$ uma função contínua ($n + 1$) vezes diferenciável num intervalo $[a, b]$. Sejam x_0, x_1, \dots, x_n ($n + 1$) pontos distintos no intervalo. Então, temos

$$a_n = f[x_0, x_1, \dots, x_n] = \sum_{i=0}^n \frac{f(x_i)}{\prod_{j=0, j \neq i}^n (x_i - x_j)}$$

Portanto, o **polinômio interpolador de Newton** pode ser escrito como:

$$p_n(x) = f[x_0] + (x - x_0)f[x_0, x_1] + (x - x_0)(x - x_1)f[x_0, x_1, x_2] + \dots + \left(\prod_{i=0}^{n-1} (x - x_i) \right) f[x_0, \dots, x_n].$$

Ou seja,

Teorema 7.5

Polinômio Interpolador de Newton. seja $f(x)$ uma função contínua e definida em x_0, x_1, \dots, x_n ($n + 1$) pontos distintos de um intervalo $[a, b]$. O polinômio de grau $\leq n$ baseado nas diferenças divididas é dado por,

$$P(x) = f[x_0] + \sum_{i=1}^n \prod_{j=0}^{i-1} (x - x_j) f[x_0, x_1, \dots, x_i]$$

interpola $f(x)$ nos pontos x_0, x_1, \dots, x_n . Sendo,



$$f[x_0, x_1, \dots, x_i] = \frac{f^{n+1}(\varepsilon)}{(n+1)!}, \varepsilon \in [a; b]$$

Desta forma, podemos aproximar uma função $f(x)$ por meio do polinômio interpolador de Newton num dado conjunto de dados definido previamente. Uma maneira mais prática e direta de obtenção do polinômio interpolador de Newton é por meio da construção de uma tabela com as diferenças divididas. Por exemplo,

x_i	$f(x_i)$	ordem um	ordem dois	ordem três	ordem quatro
x_0	$f(x_0)$	$\frac{f(x_1)-f(x_0)}{x_1-x_0}$	$\frac{f[x_1,x_2]-f[x_0,x_1]}{x_2-x_0}$		
x_1	$f(x_1)$				
x_2	$f(x_2)$	$\frac{f(x_2)-f(x_1)}{x_2-x_1}$	$\frac{f[x_2,x_3]-f[x_1,x_2]}{x_3-x_1}$	$\frac{f[x_1,x_2,x_3]-f[x_0,x_1,x_2]}{x_3-x_0}$	$\frac{f[x_1,x_2,x_3,x_4]-f[x_0,x_1,x_2,x_3]}{x_4-x_0}$
x_3	$f(x_3)$	$\frac{f(x_3)-f(x_2)}{x_3-x_2}$			
x_4	$f(x_4)$	$\frac{f(x_4)-f(x_3)}{x_4-x_3}$	$\frac{f[x_3,x_4]-f[x_2,x_3]}{x_4-x_2}$	$\frac{f[x_2,x_3,x_4]-f[x_1,x_2,x_3]}{x_4-x_1}$	

Exemplo 7.4

Construir a tabela de diferenças divididas da função,

$$f(x) = \frac{1}{x}$$

definida nos pontos $x_0 = 1, x_1 = 2, x_2 = 4$ e $x_3 = 5$.

x_i	$f(x_i)$	ordem um	ordem dois	ordem três
1	1			
2	$\frac{1}{2}$	$-\frac{1}{2}$	$\frac{1}{8}$	$-\frac{1}{40}$
4	$\frac{1}{4}$	$-\frac{1}{8}$		
5	$\frac{1}{5}$	$-\frac{1}{20}$	$\frac{1}{40}$	



Exemplo 7.5

Considere

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-x^2/2} dx.$$

Obtenha um polinômio interpolador de grau 3 para o integrando no intervalo $1, 4 \leq x \leq 2$ e veja como encontrar uma aproximação para o cálculo da integral no intervalo dado.

x_i	$f(x_i)$	ordem um	ordem dois	ordem três
1,4	0,3753			
1,6	0,2780	-0,4863	0,2142	0,00091
1,8	-0,1978	-0,4000		
2,0	0,1353	-0,3128	0,2197	

$$p_3(x) = f[x_0] + (x-x_0)f[x_0, x_1] + (x-x_0)(x-x_1)f[x_0, x_1, x_2]$$

$$+ (x-x_0)(x-x_1)(x-x_2)f[x_0, x_1, x_2, x_3]$$

$$= 0,3753 - 0,4863(x-1,4) + 0,2142(x-1,4)(x-1,6) + 0,0091(x-1,4)(x-1,6)(x-1,8).$$

Assim,

$$\frac{1}{\sqrt{2\pi}} \int_{1,4}^2 e^{-x^2/2} dx \approx \frac{1}{\sqrt{2\pi}} \int_{1,4}^2 p_3(x) dx.$$

RAs funções,

$$N_k(x) = \prod_{i=0}^k (x-x_i), \quad k = 1, \dots, n$$



formam uma base no espaço de funções polinomiais de grau n . Diferentemente da base polinomial de Lagrange, esta base e as condições $\{p_n(x_i) = f(x_i)\}_{i=0}^n$ nos leva a um sistema de equações lineares cuja matriz é triangular inferior. Ou seja,

$$\begin{pmatrix} 1 & 0 & 0 & 0 & \cdots & 0 \\ 1 & x_1 - x_0 & 0 & 0 & \cdots & 0 \\ 1 & x_2 - x_0 & (x_2 - x_0)(x_2 - x_1) & 0 & \cdots & 0 \\ \vdots & \cdots & \ddots & \vdots & & \\ 1 & x_k - x_0 & \cdots & \cdots & \prod_{j=0}^{k-1} (x_k - x_j) & \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} f(x_0) \\ f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix}$$

cuja solução é o conjunto de equações apresentadas no teorema 7.4. O método de interpolação de Newton é recomendável quando há a necessidade de atualização ou adição (remoção) de dados.

Limitante superior do erro

$$|E(x)| \leq \frac{|\prod_{i=0}^n (x - x_i)|}{(n+1)!} M$$

em que

$$M = \max |f^{(n+1)}(x)|, x \in [x_0, x_n]$$

ou

$$M = |f[x_0, x_1, \dots, x_n]|, x \in [x_0, x_n]$$



Exemplo 7.6

Considere a função $f(x) = e^x + \text{sen}(x)$, tabelada nos pontos,

x_i	0	0,5	1,0
y_i	1	2,1281	3,5598

determine o polinômio interpolador de Newton, calcule o valor da função no ponto $x = 0,7$ e o limitante superior do erro.

EXERCÍCIOS PROPOSTOS

Exercício 7.1

A tabela seguinte apresenta a velocidade de queda de um paraquedista em função do tempo:

$t(s)$	1	3	5	7	20
$v(cm/s)$	800	2310	3090	3940	8000

estime o valor da velocidade no instante de tempo $t = 10s$, utilizando um polinômio interpolador de grau 3.

Exercício 7.2

Para a função dada, seja $x_0 = 0$, $x_1 = 0,6$ e $x_2 = 0,9$. Construa polinômios de grau $n \leq 2$, para aproximar $f(0,45)$, e encontre o erro relativo.



- $f(x) = \cos(x)$
- $f(x) = \sqrt{1+x}$
- $f(x) = \ln(x+1)$

Exercício 7.3

A velocidade ascendente de um foguete é dada como uma função do tempo, conforme a seguinte tabela:

$t(s)$	0	10	15	20	22,5	30
$v(m/s)$	0	22,04	362,78	517,35	602,97	901,67

- determine a velocidade do foguete em $t = 16s$;
- determine a distância percorrida pelo foguete entre $t = 11s$ e $t = 16s$;
- determine a aceleração do foguete em $t = 16s$.

Exercício 7.4

Considere a seguinte função

$$S_i(x) = \int_0^x \frac{\text{sen}(t)}{t} dt \quad t \in [0; 10].$$

Construa um polinômio interpolador que melhor ajusta a função original. Considere $n = 6$.



Exercício 7.5

Num determinado experimento, levou-se a água ao ponto de ebulição, desejamos no entanto calcular o calor específico da água a 61°C . O calor específico da água é dado conforme a seguinte tabela:

Temperatura	22	42	52	82	100
Calor específico	4181	4179	4186	4199	4217

Determine o valor do calor específico na temperatura de 61°C utilizando os polinômios interpoladores de Lagrange.

Exercício 7.6

Os seguintes dados caracterizam despesas referentes a mão de obra e ao nível de produção

mês	mar	abr	mai	jun	jul	ago
produção (unidades)	200	300	400	640	540	580
despesas (x 1000)	2,50	2,80	3,15	3,83	3,225	3,70

Estime as despesas para uma produção de 350 unidades.



Exercício 7.7

Num concurso, o número de candidatos que obtiveram notas entre determinado intervalos foi o seguinte:

questões corretas	0 - 20	21 - 40	41 - 60	61 - 80	81 - 100	101 - 120
N. de candidatos	43	68	63	60	18	7

Estime o número de candidatos que acertaram aproximadamente 70 questões.

Exercício 7.8

A seguinte tabela fornece o ponto de fusão de uma liga metálica ao longo do tempo, caracterizando a temperatura e o percentual de um certo metal na liga.

Temperatura (°C)	180	200	250	290	300
Pressão (bar)	41	79	86	99	87

Determine o ponto de fusão da liga para um percentual de 84 por cento do metal.



UNIDADE III



INTEGRAÇÃO NUMÉRICA

A existência da integral definida de uma função f não negativa num intervalo fechado $[a; b]$ é baseada na interpretação da integral como o cálculo da área sob o gráfico de uma função no intervalo. A integral definida pode ser definida pelo conceito de integral de Riemann. Para isso, vamos particionar o intervalo $[a, b]$ de modo que a partição P seja dada por,

$$P = \{a = x_0 < x_1 < x_2 < \dots < x_{n-1} < x_n = b\}$$

que divide o intervalo $[a, b]$ em n subintervalos $[x_{i-1}, x_i]$.

Definição 8.1

Para cada subintervalo $I_i = [x_{i-1}, x_i]$ em P , seja

$$m_i = \inf\{f(x); x \in [x_{i-1}, x_i]\}$$

e

$$M_i = \sup\{f(x); x \in [x_{i-1}, x_i]\}$$

o ínfimo e o **supremo** da função

$$f|_{[x_{i-1}, x_i]}$$



Definição 8.2

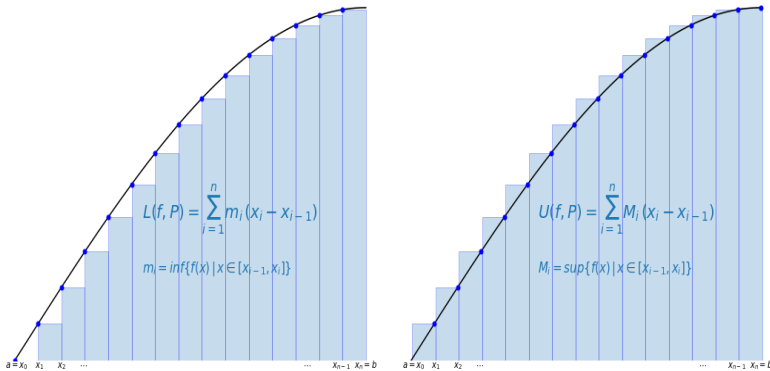
A Soma **inferior** da f com relação a partição P é dada por

$$L(f, P) = \sum_{i=1}^n m_i(x_i - x_{i-1})$$

e a soma **superior** da f com relação a P é dada por,

$$U(f, P) = \sum_{i=1}^n M_i(x_i - x_{i-1})$$

Figura 8.1 – Exemplo de uma soma inferior e superior



Claramente, temos que

$$L(f, P) \leq \underbrace{\sum_{i=1}^n f(x_i)(x_{i-1} - x_i)}_{*} \leq U(f, P)$$

em que o termo central da desigualdade (*) é denominado de soma de **Riemann**.



Teorema 8.1

Integral de Riemann.. Suponha que $f: [a, b] \rightarrow \mathbb{R}$ é contínua no intervalo, então

$$\int_a^b f(x)dx = \lim_{n \rightarrow \infty} \sum_{i=1}^n f(x_i)(x_{i-1} - x_i) \quad (8.1)$$

Portanto, pela soma de Riemann, podemos aproximar a integral definida de uma função no intervalo $[a, b]$ por meio da seguinte expressão. Ou seja,

$$\int_a^b f(x)dx \approx \sum_{i=0}^n w_i f(x_i) \quad (8.2)$$

REGRA DO TRAPÉZIO

considerando o caso mais elementar, ou seja, calcular a área sobre uma reta de modo que a função $f(x)$ pode ser aproximada via interpolação de Lagrange por um polinômio de grau 1. Ou seja,

$$\int_a^b f(x)dx \approx \int_a^b p_1(x)dx. \quad (8.3)$$



Assim,

$$\begin{aligned}\int_a^b f(x)dx &\approx \int_a^b p_1(x)dx = \int_a^b \left[f(a)\frac{(x-b)}{(a-b)} + f(b)\frac{(x-a)}{(b-a)} \right] dx \\ &= \int_a^b \left[-f(a)\frac{(x-b)}{(b-a)} + f(b)\frac{(x-a)}{(b-a)} \right] dx \\ &= \frac{f(a)}{b-a} \int_b^a (x-b)dx + \frac{f(b)}{b-a} \int_a^b (x-a)dx \\ &= \frac{f(a)}{b-a} \left[\frac{x^2}{2} - bx \right]_a^b + \frac{f(b)}{b-a} \left[\frac{x^2}{2} - ax \right]_a^b \\ &= \frac{f(a)}{b-a} \left(\frac{b^2 - 2ab + a^2}{2} \right) + \frac{f(b)}{b-a} \left(\frac{b^2 - 2ab + a^2}{2} \right) \\ &= \frac{f(a)}{b-a} \left(\frac{(b-a)^2}{2} \right) + \frac{f(b)}{b-a} \left(\frac{(b-a)^2}{2} \right).\end{aligned}$$

Ou seja,

$$\int_a^b f(x)dx \approx \frac{b-a}{2} (f(a) + f(b)) \quad (8.4)$$

que é conhecida como **regra do trapézio**. Seja,

$$E_t = \frac{-f''(\xi)h^3}{12}$$

o erro de truncamento com, $h = (b - a)$ e ξ como o valor que maximiza $|f''(\xi)|$. De maneira que podemos definir o limitante superior para o erro por,

$$|E_t| \leq \frac{h^3}{12} \max |f''(x)|$$



Exemplo 8.1

Calcule numericamente pelo método do trapézio e em seguida calcule o erro real e o erro máximo pela fórmula do erro de truncamento.

$$I_1 = \int_0^{\frac{\pi}{2}} \text{sen}(x) dx$$

$$I_1 = \frac{\frac{\pi}{2} - 0}{2} \left(\text{sen}(0) + \text{sen}\left(\frac{\pi}{2}\right) \right) = \frac{\pi}{4}$$

$$E_r = \frac{\left|1 - \frac{\pi}{4}\right|}{1} = \frac{4}{4} - \frac{\pi}{4} = \frac{4 - \pi}{4} = 0,2149$$

$$E_t = \frac{-f''(\varepsilon)h^3}{12} = \frac{\text{sen}\left(\frac{\pi}{2}\right)}{12} \left(\frac{\pi}{2} - 0\right)^3 = 0,3230$$

que corresponde ao erro máximo.

Exemplo 8.2

Para I_2 , temos:

$$I_2 = \int_3^{3,6} \frac{1}{x} dx$$

$$I_2 = \frac{3,6 - 3}{2} \left(\frac{1}{3} + \frac{1}{3,6} \right) = \frac{11}{60}$$

$$E_r = \frac{|0,1823 - 0,1833|}{0,1823} = 5,48 \times 10^{-3}$$

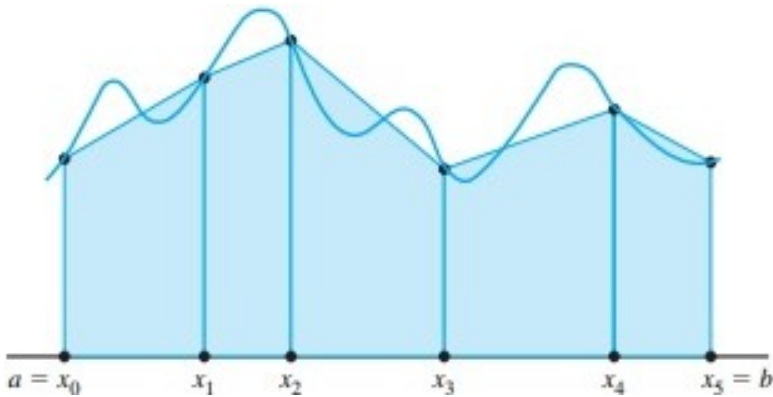
Repetindo o cálculo de I_2 , considerando o limite de integração entre 3 e 100.



$$\int_a^b f(x)dx = \sum_{i=0}^{n-1} \int_{x_i}^{x_{i+1}} f(x)dx \approx \sum_{i=0}^{n-1} \frac{h}{2} (f(x_i) + f(x_{i+1}))$$

$$\approx \sum_{i=0}^{n-1} \frac{h}{2} (f(x_0) + 2f(x_1) + 2f(x_2) + \dots + 2f(x_{n-1}) + f(x_n))$$

Figura 8.2 - Exemplo da regra do trapézio composto



que é conhecida como **regra do trapézio composto**.

$$\int_a^b f(x)dx \approx \frac{(b-a)}{2n} \left(f(a) + 2 \sum_{i=1}^{n-1} f(x_i) + f(b) \right)$$

ERRO DE TRUNCAMENTO

Em cada um dos n subintervalos, comete-se um erro ligado ao método dos trapézios. Ou seja,



$$E_t \cong n \left(\frac{-f''(\epsilon)h^3}{12} \right)$$

$$E_t \cong -\frac{b-a}{12} f''(\epsilon)h^2$$

O erro é aproximado pois ϵ varia para cada segmento. Sendo assim podemos definir o limitante superior para o erro como sendo,

$$|E_t| \leq \frac{h^2}{12}(b-a) \max |f''(\epsilon)|$$

Exemplo 8.3

Calcule numericamente a seguinte integral utilizando a regra dos trapézios composta com $n = 3$.

$$I = \int_0^1 x^2 dx$$

e, calcule numericamente a seguinte integral com o limitante superior do erro,

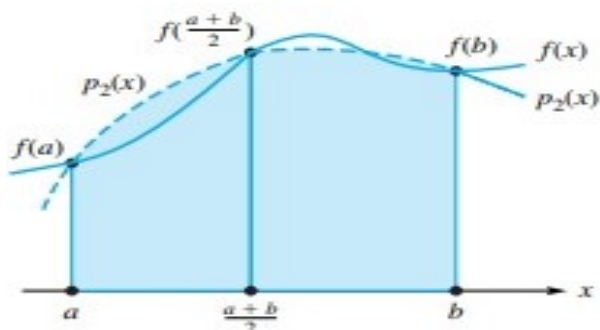
$$\int_0^1 \frac{\cos(x)}{1+x} dx$$

AS REGRAS DE SIMPSON

Considere uma função $f(x)$ definida em três pontos distintos e igualmente espaçados x_0 , x_1 e x_2 num intervalo $[a; b]$, de modo que possamos aproximar a função por um polinômio interpolador (Lagrange) de grau 2 (conforme visto na seguinte figura).



Figura 8.3 – Exemplo ilustrativo para o método de Simpson



Escrito da seguinte forma,

$$f(x) \approx p_2(x) = f(a)L_0(x) + f\left(\frac{a+b}{2}\right)L_1(x) + f(b)L_2(x)$$

tal que

$$L_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^2 \frac{(x - x_j)}{(x_i - x_j)}, \quad i = 0, 1, 2.$$

Considerando, $x_0 = a$, $x_1 = (a+b)/2$ e $x_2 = b$. Tal que,

$$\begin{cases} x_1 - x_0 = h \\ x_2 - x_1 = h \\ x_2 - x_0 = 2h. \end{cases}$$

Assim, $p_2(x)$ pode ser escrito como,



$$\begin{aligned}
 p_2(x) &= f(x_0) \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} + f(x_1) \frac{(x-x_0)(x-x_1)}{(x_1-x_0)(x_1-x_2)} + f(x_2) \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)} \\
 &= f(x_0) \frac{(x-x_1)(x-x_2)}{(-h)(-2h)} + f(x_1) \frac{(x-x_0)(x-x_1)}{h(-h)} + f(x_2) \frac{(x-x_0)(x-x_1)}{2h(h)} \\
 &= \frac{f(x_0)}{2h^2}(x-x_1)(x-x_2) - \frac{f(x_1)}{h^2}(x-x_0)(x-x_2) + \frac{f(x_2)}{2h^2}(x-x_0)(x-x_1).
 \end{aligned}$$

Portanto,

$$\begin{aligned}
 \int_a^b f(x)dx &= \int_{x_0}^{x_2} f(x)dx \approx \int_{x_0}^{x_2} p_2(x)dx \\
 &= \int_{x_0}^{x_2} p_2(x)dx =
 \end{aligned}$$

onde,

$$\frac{f(x_0)}{2h^2} \int_{x_0}^{x_2} (x-x_1)(x-x_2)dx - \frac{f(x_1)}{h^2} \int_{x_0}^{x_2} (x-x_0)(x-x_2)dx + \frac{f(x_2)}{2h^2} \int_{x_0}^{x_2} (x-x_0)(x-x_1)dx.$$

De modo que,

$$\int_{x_0}^{x_2} f(x)dx \approx \int_{x_0}^{x_2} p_2(x)dx = \frac{h}{3} (f(x_0) + 4f(x_1) + f(x_2))$$

ou,

$$\int_a^b f(x)dx \approx \frac{h}{3} \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right)$$

que é denominada de **primeira regra de Simpson** ou **regra do 1/3**. Onde $h = (b-a)/2$.



LIMITANTE SUPERIOR PARA O ERRO

$$|E_2| \leq \frac{h^5}{90} \max |f^4(\varepsilon)|$$

REGRA 1/3 DE SIMPSON GENERALIZADA

Consiste em dividirmos o intervalo $[a; b]$ em n subintervalos de tamanho constante (h) e a cada par de subintervalos aplicar a primeira regra de Simpson. **OBS:** O número n de subintervalos deverá sempre ser par, de maneira que:

$$\int_a^b f(x) dx \approx$$

$$\frac{h}{3} [f(x_0) + 4f(x_1) + f(x_2)] + \frac{h}{3} [f(x_2) + 4f(x_3) + f(x_4)] + \dots +$$
$$\frac{h}{3} [f(x_{n-2}) + 4f(x_{n-1}) + f(x_n)]$$

$$\int_a^b f(x) dx \approx$$

$$\frac{h}{3} [f(x_0) + 4f(x_1) + f(x_2)] + \frac{h}{3} [f(x_2) + 4f(x_3) + f(x_4)] + \dots +$$
$$\frac{h}{3} [f(x_{n-2}) + 4f(x_{n-1}) + f(x_n)]$$

$$\int_a^b f(x) dx \approx \frac{h}{3} \sum_{i=0}^n [f(x_{2i}) + 4f(x_{2i+1}) + f(x_{2i+2})]$$

com,

$$h = \frac{(b-a)}{n}$$



LIMITANTE SUPERIOR PARA O ERRO

$$|E_t| \leq \frac{h^4}{180} (x_n - x_0) \max |f^4(\xi)|$$

Exemplo 8.4

Calcule numericamente a integral:

$$I = 4 \int_0^1 \frac{dx}{1+x^2}$$

utilizando a primeira regra de Simpson para $n = 2$ e para $n = 6$.

Calcule o valor de π , dado pela expressão:

$$I = 4 \int_0^1 \frac{dx}{1+x^2}$$

com um erro menor que 10^{-4} .

Exemplo 8.5

Calcular o trabalho realizado por um gás sendo aquecido segundo a tabela,

$V(m^3)$	1,5	2,0	2,5	3,0	3,5	4,0	4,5
$P(\frac{kg}{m^2})$	80	72	64	53	44	31	22

onde,

...

$$W = \int_{V_i}^{V_f} P dV$$



REGRA 3/8 OU SEGUNDA REGRA DE SIMPSON

A regra 3/8 de Simpson consiste em aproximar a função $f(x)$ entre quatro pontos consecutivos x_0, x_1, x_2 e x_3 por um polinômio de grau 3.

Neste caso, teremos a seguinte aproximação:

$$\int_a^b f(x)dx \approx \int_{x_0}^{x_3} p_3(x)dx = \frac{3h}{8} [f(x_0) + 3f(x_1) + 3f(x_2) + f(x_3)]$$

$$h = \frac{(b-a)}{n}.$$

O erro de truncamento é:

$$E_t = -\frac{3f^{(4)}(\xi)h^5}{80}$$

Considere a divisão do intervalo $[a, b]$ num número múltiplo de 3 ($3n$) subintervalos de tamanho $h = \frac{(b-a)}{n}$.

Neste caso, teremos a seguinte aproximação:

$$\begin{aligned} \int_a^b f(x)dx &= \int_{x_0}^{x_n} f(x)dx \\ &\approx \frac{3h}{8} [f(x_0) + 3f(x_1) + 3f(x_2) + f(x_3)] \\ &+ \frac{3h}{8} [f(x_3) + 3f(x_4) + 3f(x_5) + f(x_6)] + \dots \\ &+ \frac{3h}{8} [f(x_{n-3}) + 3f(x_{n-2}) + 3f(x_{n-1}) + f(x_n)] \end{aligned}$$



QUADRATURA DE GAUSS

Considere a integral

$$\int_a^b f(x) dx$$

desejamos desenvolver uma fórmula de integração da forma:

$$\int_a^b f(x) dx = w_0 f(x_0) + w_1 f(x_1) + \cdots + w_n f(x_n).$$

se $f(x) = p_m(x)$,

$$\int_a^b p_m(x) dx \approx \sum_{i=0}^n w_i p_m(x_i)$$

IDEIA PRINCIPAL

Os pontos $a = x_0 < x_1 < \cdots < x_n = b$ não são necessariamente equidistantes,

$$\int_a^b f(x) dx \cong w_0 f(x_0) + w_1 f(x_1) + \cdots + w_n f(x_n), \forall f \in p_n$$

escolher um polinômio de maior grau possível,

$$p_m(x) = a_m x^m + a_{m-1} x^{m-1} + \cdots + a_1 x + a_0$$

de modo que: $f(x) = 1, x, x^2, x^3, \dots, x^m \in p_m$,

com a escolha apropriada de valores e constantes, podemos obter que a aproximação seja exata nesse conjunto. A Quadratura de Gauss mais comum é a de Gauss - Legendre, por estarem associadas ao polinômio de Legendre.



$$P_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} [(x^2 - 1)^n].$$

Os polinômios de Legendre são ortogonais com relação ao produto interno definido por,

$$\langle f, g \rangle = \int_{-1}^1 f(x)g(x)dx$$

os quatro primeiros polinômios de Legendre são:

$$p_0(x) = 1, p_1(x) = x, p_2(x) = x^2 - \frac{1}{3}, p_3(x) = x^3 - \frac{3}{5}x.$$

PROPRIEDADES BÁSICAS DOS POLINÔMIOS DE LEGENDRE

- Os polinômios de Legendre são ortogonais a outros polinômios. Ou seja,

$$\langle p_n(x), q_m(x) \rangle = \int_{-1}^1 p_n(x)q_m(x)dx = 0, n > m,$$

sendo $q_m(x)$ um polinômio qualquer de grau menor que n .

- Os polinômios de Legendre são ortogonais entre si,

$$\langle p_n(x), p_m(x) \rangle = \int_{-1}^1 p_n(x)p_m(x)dx = 0, n \neq m.$$

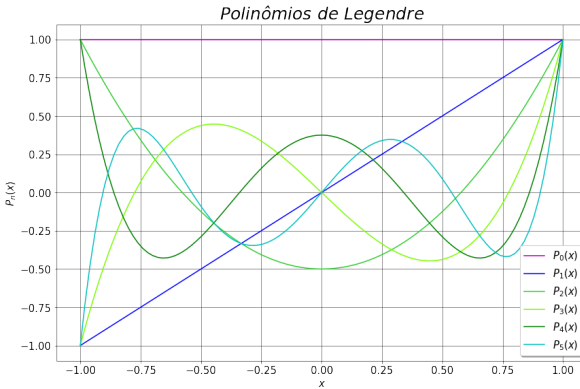
- Se os polinômios de Legendre forem iguais, então

$$\langle p_n(x), p_n(x) \rangle = \int_{-1}^1 \|p_n(x)\|^2 dx = \frac{2}{2n+1}.$$



- $p_n(1) = 1$ e $p_n(-1) = (-1)^n$, $n = 0, 1, 2, \dots$.
- O polinômio de Legendre $p_n(x)$ de grau $n \geq 1$ possui n zeros reais, distintos, pertencentes ao intervalo $(-1, 1)$ e simétricos em relação à origem.

Figura 8.4 - Polinômios de Legendre da ordem zero a ordem 5



QUADRATURA DE GAUSS PARA $N = 1$

$$\int_a^b f(x) \cong w_0 f(x_0), \quad a \leq x_0 \leq b$$

$$f(x) \approx p_1(x) = a_0 + a_1 x$$

$$\int_a^b (a_0 + a_1 x) dx = a_0 x \Big|_a^b + a_1 \frac{x^2}{2} \Big|_a^b = a_0(b-a) + a_1 \frac{(b^2 - a^2)}{2} \quad (8.6)$$

e,

$$w_0 f(x_0) = w_0(a_0 + a_1 x_0) = a_0 w_0 + a_1 w_0 x_0 \quad (8.7)$$



igualando as equações (8.6) com a (8.7) obtemos:

$$w_0 = b - a \text{ e } x_0 = \frac{b+a}{2}$$

de modo que,

$$\int_a^b f(x)dx \cong (b-a)f\left(\frac{b+a}{2}\right) \iff \frac{b-a}{2}(f(a) + f(b))$$

que corresponde simplesmente a **regra do trapézio** já estudada.

QUADRATURA DE GAUSS PARA $N = 2$ E $N = 3$

$$\int_a^b f(x) \cong w_0 f(x_0) + w_1 f(x_1)$$

$$f(x) \approx p_2(x) = a_0 + a_1 x + a_2 x^2$$

$$\int_a^b (a_0 + a_1 x + a_2 x^2) dx = a_0 x \Big|_a^b + a_1 \frac{x^2}{2} \Big|_a^b + a_2 \frac{x^3}{3} \Big|_a^b =$$

$$a_0(b-a) + a_1 \frac{(b^2 - a^2)}{2} + a_2 \frac{(b^3 - a^3)}{3}$$

se considerarmos a seguinte transformação,

$$\int_a^b f(x) dx = \int_{-1}^1 f(x) dx = 2a_0 + \frac{a_1}{2}(1^2 - (-1)^2) + \frac{a_2}{3}(1^3 - (-1)^3) = 2a_0 + \frac{2}{3}a_2.$$

Com isso,

$$w_0 f(x_0) + w_1 f(x_1) = w_0 (a_0 + a_1 x_0 + a_2 x_0^2) + w_1 (a_0 + a_1 x_1 + a_2 x_1^2)$$

$$= a_0(w_0 + w_1) + a_1(w_0 x_0 + w_1 x_1) + a_2(w_0 x_0^2 + w_1 x_1^2).$$



De modo que, das equações anteriores, podemos tirar as seguintes relações,

$$\begin{cases} w_0 + w_1 = 2 \\ w_0 x_0 + w_1 x_1 = 0 \\ w_0 x_0^2 + w_1 x_1^2 = 2/3 \end{cases}$$

Considerando $w_0 = w_1 = 1$. Temos que, $x_0 + x_1 = 0 \rightarrow x_0 = -x_1$ e com isto, podemos inferir através da última relação que,

$$x_0^2 + x_1^2 = \frac{2}{3} \rightarrow (-x_1)^2 + x_1^2 = \frac{2}{3} \rightarrow x_1 = \pm \frac{1}{\sqrt{3}} = \pm \frac{\sqrt{3}}{3}.$$

Assim, concluímos que:

$$\int_{-1}^1 f(x) dx \cong f\left(-\frac{1}{\sqrt{3}}\right) + f\left(\frac{1}{\sqrt{3}}\right)$$

de maneira similar para $n = 3$, temos:

$$\int_{-1}^1 f(x) dx \cong \frac{1}{9} \left(5f\left(-\sqrt{0,6}\right) + 8f(0) + 5f\left(\sqrt{0,6}\right) \right)$$

uma integral

$$\int_a^b f(x) dx, [a, b]$$

pode ser transformada em uma integral em $[-1, 1]$ utilizando mudança de variável. Ou seja,

$$t = \frac{2x - (a+b)}{b-a} \iff x = \frac{1}{2} [(b-a)t + a + b].$$



Facilmente visto que se $t = -1 \rightarrow x = a$ e para $t = 1 \rightarrow x = b$. Portanto, a mudança de variável é possível. Com isso a quadratura de Gauss pode ser aplicada a todo intervalo $[a, b]$. Assim, aplicando a seguinte transformação:

$$t = \frac{b-a}{2}x + \frac{b+a}{2}, \quad -1 \leq x \leq 1.$$

Onde, $dt = \frac{b-a}{2} dx$, então

$$\int_a^b f(t) dt = \frac{b-a}{2} \int_{-1}^1 f\left(\frac{b-a}{2}x + \frac{b+a}{2}\right) dx$$

f^x ,

$$\int_a^b f(x) dx = \frac{b-a}{2} \int_{-1}^1 f\left(\frac{b-a}{2}t + \frac{b+a}{2}\right) dt.$$

EXERCÍCIOS PROPOSTOS

Exercício 8.1

Calcule a seguinte integral

$$I = \int_0^{\pi} e^{\cos x} dx$$

com erro menor ou igual a 10^{-3} .



Exercício 8.2

Calcule a seguinte integral

$$I = \int_0^{2\pi} f(x) dx,$$

onde,

$$f(x) = \begin{cases} \text{sen}(x), & 0 \leq x \leq \frac{\pi}{2}, \\ \text{cos}(x), & \frac{\pi}{2} \leq x \leq 2\pi. \end{cases}$$

Exercício 8.3

Calcule a seguinte integral

$$I = \int_0^2 e^{2x} \text{sen}(6x) dx$$

Exercício 8.4

Um Boeing 727-200 de massa $m = 97000$ kg aterrissa a uma velocidade de 93 m/s (em torno de 181 nós) e liga os seus reversos em $t = 0$. A força F aplicada no avião à medida que ele reduz a sua velocidade é dada por $F = -(5v^2 + 570000)$, onde v é a velocidade do avião. Usando a segunda lei de Newton do movimento e da dinâmica dos fluidos,



a relação entre a velocidade e a posição x do avião pode ser escrita como:

$$mv \frac{dv}{dx} = F$$

onde x é a distância medida a partir da localização do jato em $t = 0$. Determine a distância percorrida pelo avião antes que sua velocidade se reduza a 40 m/s (em torno de 78 nós) usando o método trapezoidal composto.

Exercício 8.5

A função $f(x)$ é dada na forma tabulada a seguir.

x	0	0,25	0,5	0,75	1,0
$f(x)$	0,9162	0,8109	0,6931	0,5596	0,4055

compare

$$\int_0^1 f(x) dx \text{ com } h=0,25 \text{ e } h=0,5.$$

Utilizando,

O método trapezoidal composto.



- Use interpolação linear para determinar $f(x)$ nos pontos centrais.
- A regra de Simpson 1/3 generalizada.

Exercício 8.6

Para estimar a área superficial de uma bola de futebol americano, mede-se o seu diâmetro em diferentes pontos. A área superficial S e o volume V podem ser determinados usando:

$$S = 2\pi \int_0^L r dz$$

e

$$V = \pi \int_0^L r^2 dz$$

Use os dados abaixo para determinar o volume e a área superficial da bola.

$z(\text{cm})$	0	1	2	3	4	5	6	7	8	9	10	11	12
$d(\text{cm})$	0	13	16	24	28	30	31	30	28	24	17	13	0



Exercício 8.7

O perímetro P de uma elipse é dado por:

$$P = 4a \int_0^{\frac{\pi}{2}} \sqrt{1 - k^2 \sin^2(\theta)} d\theta$$

onde,

$$k = \frac{\sqrt{a^2 + b^2}}{a}.$$

Com a e b sendo, aos eixos maior e menor, respectivamente. Calcule o perímetro da seguinte elipse:

$$\frac{x^2}{5^2} + \frac{y^2}{4^2} = 1.$$

Exercício 8.8

A densidade ρ da terra varia com o raio r . A seguinte tabela fornece a densidade aproximada em diferentes raios:

$r(km)$	0	800	1200	1400	2000	3000	3400	3600	4000
$\rho(\frac{kg}{m^3})$	13000	12900	12700	12000	11650	10600	9900	5500	5300
$r(km)$	5000	5500	6370						
$\rho(\frac{kg}{m^3})$	4750	4500	3300						



A massa da terra pode ser calculada por meio da seguinte integral:

$$m = 4\pi \int_0^{6370} \rho r^2 dr.$$

Escreva uma função que calcule a massa da terra a partir dos dados tabelados com um espaçamento qualquer.

Exercício 8.9

A seguinte função exibe tantos regiões achatadas quanto íngremes em uma região relativamente curva de x ,

$$f(x) = \frac{1}{(x-0,3)^2 + 0,01} + \frac{1}{(x-0,9)^2 + 0,04} - 6.$$

Determine o valor da integral definida dessa função com x variando de 0 a 1 e $h = 0,25$.

Exercício 8.10

A forma da linha centróide de um dado arco de comprimento L , pode ser modelada aproximadamente pela equação:

$$f(x) = 211,5 - 20,97 \cosh\left(\frac{x}{30,38}\right), \quad -91,21 \leq x \leq 91,21.$$



Sabendo que,

$$L = \int_a^b \sqrt{1 + \left[\frac{df}{dx}(x) \right]^2} dx.$$

Determine o comprimento do arco.

Exercício 8.11

Num dado sistema de comunicação com sinalização multi-dimensional (ortogonal), o BER (Bit Error Rate) é derivado a partir da seguinte função distribuição de probabilidade,

$$P_{e,b} = \frac{2^{b-1}}{2^b - 1} \left(1 - \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} \left(Q^{M-1} \left(-\sqrt{2}y - \sqrt{bSNR} \right) \right) e^{-y^2} dy \right)$$

onde, b representa o número de bits, $M = 2^b$ a quantidade de sinais ortogonais, SNR a relação sinal - ruído, e $Q(\cdot)$ é a função erro definida por,

$$Q(x) = \frac{1}{\sqrt{2\pi}} \int_x^{\infty} e^{-y^2/2}.$$

Obtenha uma aproximação do valor de $P_{e,b}(SNR, b)$ para uma relação sinal - ruído de 8 dB e 4 bits.



Exercício 8.12

O Volume de uma estrutura formada por uma superfície de revolução de uma função $y = f(x)$ em torno do eixo x num intervalo fechado $[a, b]$, pode ser descrita por meio da seguinte integral,

$$I = \pi \int_a^b f^2(x) dx.$$

Calcule o volume de uma esfera com raio de comprimento unitário.

Exercício 8.13

Uma linha reta foi traçada de modo a tangenciar as margens de um rio nos pontos A e B . Para medir a área do trecho entre o rio e a reta AB foram traçadas perpendiculares em relação a AB com intervalos de $0,05$ m. Determine a área aproximada no trecho. Sabendo que,

<i>perpendiculares</i>	1	2	3	4	5	6	7	8
<i>comprimento (m)</i>	3,28	4,02	4,64	5,26	4,98	3,62	3,82	4,68
	9	10	11					
	5,26	3,82	3,24					



RESOLUÇÃO NUMÉRICA DE EDOS

Uma EDO (Equação Diferencial Ordinária) é uma equação regida por uma ou mais derivadas de uma função. São geralmente utilizadas para modelar e descrever a dinâmica em inúmeros modelos nas Ciências e na Engenharia.

Em alguns casos, tais equações não apresentam uma solução exata, sendo necessário a utilização de algum método numérico de modo a apresentar uma solução aproximada ao problema, que dentro de uma certa margem de erro, pode ser bastante útil na análise e solução de tais problemas. Se definirmos um valor inicial na análise de tais equações, teremos então um problema de PVI (Problema de Valor Inicial) que irá nos auxiliar na resolução da EDO, que consiste em encontrar uma função $y(t)$ que satisfaz,

$$y' = \frac{dy(t)}{dt} = f(t, y(t)) \quad (9.1)$$

com a condição inicial $y(t_0) = y_0$.

Numa solução numérica o objetivo se resume a encontrar uma sequência em que a solução analítica é aproximada por uma solução numérica convergente. Ou seja,

$$\|y_i - y(t_i)\| \leq tol$$

com tamanho de passo $h = t_{i+1} - t_i$ e o erro controlado por uma tolerância definida.



PROBLEMA DE VALOR INICIAL - PVI

A equação diferencial ordinária (9.1) sem uma condição inicial geralmente possui uma família de soluções. Ao especificar uma condição inicial, podemos identificar qual função dentre as apresentadas na famílias estamos interessados. Um PVI para uma equação diferencial ordinária de primeira ordem é a equação juntamente com uma condição inicial em um intervalo específico $a \leq t \leq b$, em que,

$$y' = f(t, y(t)), \quad y(t_0) = y_0.$$

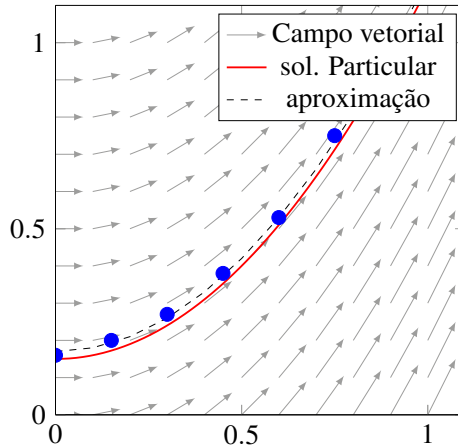
onde y' define a inclinação da reta tangente a uma curva $y(t)$ num ponto (t, y) . Podemos usar essa informação para conhecer a evolução ou a forma da solução $y(t)$. Por exemplo, a equação

$$y' = 2y$$

nos diz que a inclinação da reta tangente num ponto (t, y) é igual a sua coordenada y . Podemos também pensar a solução de uma equação diferencial como um campo vetorial de isóclinas, como pode ser vista na seguinte figura, e a equação (9.1) pode ser vista como uma inclinação para quaisquer valores atuais de (t, y) . Se usarmos um vetor para plotar a inclinação em cada ponto do plano, obtemos o campo de inclinação ou o campo direcional, da equação diferencial.



Figura 9.1 – Exemplo ilustrativo de uma família de soluções e uma solução particular com aproximação para a equação diferencial $y' = 2y$



A geometria da figura (9.1) sugere uma abordagem alternativa para “resolver” computacionalmente a equação diferencial por meio de vetores. Começando de uma condição inicial (t_0, y_0) e seguindo iterativamente na direção especificada. Uma vez efetuado o cálculo, reavalia-se a inclinação no novo ponto (t_1, y_1) e assim sucessivamente, de modo a obtermos uma boa aproximação à solução do problema do valor inicial.

MÉTODO DE EULER

O objetivo do método de Euler é obter uma aproximação para a solução do PVI

$$y' = f(t, y(t)), a \leq t \leq b, y(t_0) = y_0$$



Considerando t_0 um ponto de referência e h um número positivo. De modo que,

$$t_i = t_0 \pm ih, \quad i = 1, 2, \dots$$

as aproximações da função $y(t)$ serão calculadas nesses pontos. Se assumirmos que a função $y(t)$ possui derivadas de ordem $n + 1$ em t , sua expansão em série de Taylor é dada por,

$$y(t+h) = y(t) + hy'(t) + \frac{h^2}{2!}y''(t) + \dots + \frac{h^n}{n!}y^{(n)}(t) + \frac{h^{n+1}}{(n+1)!}y^{(n+1)}(\varepsilon)$$

onde $t < \varepsilon < t + h$.

O último termo da série representa o erro da aproximação de $y(t+h)$ pelos $n + 1$ termos da série.

Truncando a série de Taylor no termo de ordem dois, obtemos:

$$y(t+h) = y(t) + hy'(t) + \frac{h^2}{2!}y''(\varepsilon)$$

pela equação (9.1) $y' = f(t, y(t))$, temos:

$$y(t+h) = y(t) + hf(t, y(t)) + \frac{h^2}{2!}y''(\varepsilon)$$

como $h = t_{i+1} - t_i$ substituindo na equação anterior, obtemos:

$$y(t_i + t_{i+1} - t_i) = y(t_i) + (t_{i+1} - t_i)f(t_i, y(t_i)) + \frac{(t_{i+1} - t_i)^2}{2!}y''(\varepsilon)$$

com isso concluímos que,

$$y(t_{i+1}) = y(t_i) + hf(t_i, y(t_i)) \quad (9.2)$$



corresponde ao método de **Euler** para resolução numérica de EDO's de primeira ordem com PVI.

Exemplo 9.1

Resolver o PVI

$$\begin{cases} y' = y - t^2 + 1; \\ y(0) = 0,5 \\ 0 \leq t \leq 1 \\ h = 0,2. \end{cases}$$

i	t_i	$y(t_i)$	$y(t_{i+1})$
0	0,0	0,5	0,8
1	0,2	0,8	1,152
2	0,4	1,152	1,5504
3	0,6	1,5504	1,9885
4	0,8	1,9885	2,4582
5	1,0	2,4582	2,9498

$$y_1 = y_0 + h(y_0 - t_0^2 + 1) = 0,5 + 0,2(0,5 - 0 + 1) = 0,8$$

$$y_2 = y_1 + h(y_1 - t_1^2 + 1) = 0,8 + 0,2(0,8 - 0,2^2 + 1) = 1,152$$

$$y_3 = y_2 + h(y_2 - t_2^2 + 1) = 1,152 + 0,2(1,152 - 0,4^2 + 1) = 1,5504$$

$$y_4 = y_3 + h(y_3 - t_3^2 + 1) = 1,5504 + 0,2(1,5504 - 0,6^2 + 1) = 1,9885$$

$$y_5 = y_4 + h(y_4 - t_4^2 + 1) = 1,9885 + 0,2(1,9885 - 0,8^2 + 1) = 2,4582$$

$$y_6 = y_5 + h(y_5 - t_5^2 + 1) = 2,4582 + 0,2(2,4582 - 1^2 + 1) = 2,9498$$



ESTIMATIVA DO ERRO PARA O MÉTODO DE EULER

Iremos utilizar a expansão em Série de Taylor para obtermos uma estimativa do erro obtido na aproximação utilizando o método de Euler. Ou seja,

$$E_t = \sum_{k=1}^n \frac{f^{(k)}(t_k, y_k)}{(k+1)!} h^{(k+1)}$$

que no método de Euler, se resume a

$$E_t = \frac{f'(t_k, y_k)}{2} h^2$$

MÉTODO DE EULER MODIFICADO

Considerando a equação (9.1) e integrando ambos os lados da mesma, temos:

$$\int_{t_i}^{t_{i+1}} y' dt = \int_{t_i}^{t_{i+1}} f(t_i, y(t_i)) dt$$

aplicando a regra do trapezio, obtemos:

$$y(t_{i+1}) - y(t_i) = \frac{t_{i+1} - t_i}{2} (f(t_i, y(t_i)) + f(t_{i+1}, y(t_{i+1})))$$

$$\begin{aligned} y(t_{i+1}) &= y(t_i) + \frac{h}{2} (f(t_i, y(t_i)) + f(t_{i+1}, y(t_{i+1}))) \\ &= y(t_i) + \frac{h}{2} (f(t_i, y(t_i)) + f((t_i + h), y(t_i) + h \cdot f(t_i, y(t_i)))) \end{aligned}$$

considerando as variáveis auxiliares k_1 e k_2 . Tais que,



$$\begin{cases} k_1 = f(t_i, y(t_i)) \\ k_2 = f(t_i + h, y(t_i) + h.k_1) \\ y(t_{i+1}) = y(t_i) + \frac{h}{2} (k_1 + k_2) \end{cases} \quad (9.3)$$

correspondem ao método de **Euler modificado**.

Exemplo 9.2

Usando o método de Euler modificado, calcule uma aproximação para o PVI,

$$\begin{cases} y' = x - y + 2 \\ y(0) = 2 \\ x \in [0; 1], h = 0,2 \end{cases}$$

Tabela 9.1 – Tabela solução do exemplo 9.2

i	x_i	k_1	k_2	$y(x_{i+1})$
0	0,0	0,0	0,2	2,02
1	0,2	0,18	0,344	2,0724
2	0,4	0,3276	0,4621	2,1514
3	0,6	0,4486	0,5588	2,2521
4	0,8	0,5478	0,6382	2,3708

Para $i = 0$,

$$k_1 = f(x_i, y_i) = f(x_0, y_0) = x_0 - y_0 + 2 = 0 - 2 + 2 = 0$$

$$k_2 = f(x_i + h, y_i + h.k_1) = (x_0 + h) - (y_0 + h.k_1) + 2 = 0,2 - (2 + 0,2 \cdot 0) + 2 = 0,2$$

$$y(t_{i+1}) = y(t_i) + \frac{h}{2} (k_1 + k_2) = 2 + \frac{0,2}{2} (0 + 0,2) = 2,02.$$



Para $i = 1$,

$$k_1 = f(x_i, y_i) = f(x_1, y_1) = x_1 - y_1 + 2 = 0,2 - 2,02 + 2 = 0,18$$

$$k_2 = f(x_i + h, y_i + h.k_1) = (x_1 + h) - (y_1 + h.k_1) + 2 = 0,4 - (2,02 + 0,2.0,18) + 2 = 0,344$$

$$y(t_{i+1}) = y(t_i) + \frac{h}{2}(k_1 + k_2) = 2,02 + \frac{0,2}{2}(0,18 + 0,344) = 2,0724.$$

Para $i = 2$,

$$k_1 = f(x_i, y_i) = f(x_2, y_2) = x_2 - y_2 + 2 = 0,4 - 2,0724 + 2 = 0,3276$$

$$k_2 = f(x_i + h, y_i + h.k_1) = (x_2 + h) - (y_2 + h.k_1) + 2 = 0,6 - (2,0724 + 0,2.0,3276) + 2 = 0,4621$$

$$y(t_{i+1}) = y(t_i) + \frac{h}{2}(k_1 + k_2) = 2,0724 + \frac{0,2}{2}(0,3276 + 0,4621) = 2,1514.$$

Para $i = 3$,

$$k_1 = f(x_i, y_i) = f(x_3, y_3) = x_3 - y_3 + 2 = 0,6 - 2,1514 + 2 = 0,4486$$

$$k_2 = f(x_i + h, y_i + h.k_1) = (x_3 + h) - (y_3 + h.k_1) + 2 = 0,8 - (2,1514 + 0,2.0,4486) + 2 = 0,5588$$

$$y(t_{i+1}) = y(t_i) + \frac{h}{2}(k_1 + k_2) = 2,1514 + \frac{0,2}{2}(0,4486 + 0,5588) = 2,2521.$$

Como a Equação (9.3) foi obtida diretamente da regra do trapézio, o erro de truncamento local é dado por por,

$$E_t = -\frac{f''(\xi)}{12}h^3.$$

em que $x_i \leq \xi \leq x_{i+1}$. Além disso, os erros local e global são $O(h^3)$ e $O(h^2)$, respectivamente. Portanto, a diminuição do tamanho do passo h faz o erro decrescer mais rapidamente que no método de Euler.



MÉTODOS DE RUNGE - KUTTA

Os métodos de Runge–Kutta são uma família de métodos de solução de EDO's que incluem os métodos de Euler e o Euler modificado, bem como outros métodos sofisticados de ordem superior. Existem inúmeras variações na literatura, mas todas podem ser escritas por meio da forma geral:

$$\begin{cases} y(t_{i+1}) = y(t_i) + h \cdot \varphi(t_i, y(t_i), h) \\ \varphi(t_i, y(t_i), h) = \alpha_1 k_1 + \alpha_2 k_2 + \dots + \alpha_n k_n \\ \alpha_1 k_1 + \alpha_2 k_2 + \dots + \alpha_n k_n = 1 \end{cases}$$

com,

$$\begin{cases} k_1 = f(t_i, y(t_i)) \\ k_2 = f(t_i + p_1 h, y(t_i) + h(q_{11} k_1)) \\ k_3 = f(t_i + p_2 h, y(t_i) + h(q_{21} k_1 + q_{22} k_2)) \\ k_4 = f(t_i + p_3 h, y(t_i) + h(q_{31} k_1 + q_{32} k_2 + q_{33} k_3)) \\ \vdots \\ k_n = f(t_i + p_n h, y(t_i) + h(q_{n-1,1} k_1 + \dots + q_{n-1,n-1} k_{n-1})) \end{cases}$$

em que os p 's e q 's são constantes e como cada k é um cálculo da função, essa recorrência torna os métodos Runge - Kutta (ou métodos RK) eficientes na solução numérica de EDO's.



MÉTODOS DE RUNGE-KUTTA DE SEGUNDA ORDEM - RK2

O método de Runge-Kutta de Segunda Ordem é dado por:

$$y(t_{i+1}) = y(t_i) + h(\alpha_1 k_1 + \alpha_2 k_2) \quad (9.4)$$

com

$$\begin{cases} \alpha_1 + \alpha_2 = 1 \\ k_1 = f(t_i, y(t_i)) \\ k_2 = f(t_i + p_1 h, y(t_i) + h(q_{11} k_1)) \quad p_1 = q_{11} \end{cases}$$

para determinar as constantes α_1 , α_2 e p_1 , iremos desenvolver k_2 via série de Taylor, em torno do ponto de operação $(t_i, y(t_i))$ até o termo de segunda ordem, ou seja,

$$y(t_{i+1}) = y(t_i) + f(t_i, y(t_i))h + \frac{f'(t_i, y(t_i))}{2}h^2 \quad (9.5)$$

mas, pela regra da cadeia sabemos que

$$f'(t_i, y(t_i)) = f_x(t_i, y(t_i)) + f_y(t_i, y(t_i))y' \quad (9.6)$$

substituindo a expressão (9.6) na expressão (9.5), obtemos:

$$y(t_{i+1}) = y(t_i) + f(t_i, y(t_i))h + \frac{f_x(t_i, y(t_i)) + f_y(t_i, y(t_i))y'}{2}h^2 \quad (9.7)$$



como a expansão da série de Taylor de uma função em duas variáveis, é dada por:

$$f(t_i + h, y(t_i) + k) = f(t_i, y(t_i)) + (f_x(t_i, y(t_i))h + f_y(t_i, y(t_i))k) + \frac{1}{2} (f_{xx}(t_i, y(t_i))h^2 + 2f_{xy}(t_i, y(t_i))hk + f_{yy}(t_i, y(t_i))k^2) + \dots$$

e k_2 expandida via série de Taylor é dada por,

$$k_2 = f(t_i + p_1 h, y(t_i) + q_{11} k_1 h) = f(t_i, y(t_i)) + p_1 h f_x(t_i, y(t_i)) + q_{11} k_1 h + \mathcal{O}(h^2) \quad (9.8)$$

que juntamente com a expressão de k_1 , podem ser substituídas na expressão (9.4), de modo a termos a seguinte equação,

$$y(t_{i+1}) = y(t_i) + [\alpha_1 f(t_i, y(t_i)) + \alpha_2 f(t_i, y(t_i))]h + [\alpha_2 p_1 f_x + \alpha_2 q_{11} f_x f(t_i, y(t_i))]h^2 + \mathcal{O}(h^3) \quad (9.9)$$

comparando termo a termo nas equações (9.9) com (9.7), podemos tirar a seguinte relação:

$$\begin{cases} \alpha_1 + \alpha_2 = 1 \\ \alpha_2 p_1 = \frac{1}{2} \\ \alpha_2 q_{11} = \frac{1}{2} \end{cases}$$

com o sistema de equações apresenta mais equações que incógnitas, podemos obter inúmeras soluções, dentre elas podemos escolher,

$$\alpha_2 = 1 \rightarrow \alpha_1 = 0 \text{ e } p_1 = q_{11} = \frac{1}{2}.$$

Obtendo,

$$y(t_{i+1}) = y(t_i) + hk_2 \quad (9.10)$$



com

$$\begin{cases} k_1 = f(t_i, y(t_i)) \\ k_2 = f\left(t_i + \frac{h}{2}, y(t_i) + \frac{h}{2}k_1\right) \end{cases} \quad (9.11)$$

correspondendo ao método do **ponto médio**. considerando,

$$\alpha_2 = \frac{2}{3} \rightarrow \alpha_1 = \frac{1}{3} \text{ e } p_1 = q_{11} = \frac{3}{4}.$$

Teremos, então:

Definição 9.1

$$y(t_{i+1}) = y(t_i) + \frac{h}{3}(k_1 + 2k_2) \quad (9.12)$$

com

$$\begin{cases} k_1 = f(t_i, y(t_i)) \\ k_2 = f\left(t_i + h\frac{3}{4}, y(t_i) + h\frac{3}{4}k_1\right) \end{cases} \quad (9.13)$$

que corresponde ao método de **Ralston**.



Exemplo 9.3

Usando o método do ponto médio e o de Ralston resolva o seguinte PVI,

$$\begin{cases} y' = x - y + 2 \\ y(0) = 2 \\ t \in [0; 1], h = 0,2 \end{cases}$$

Tabela 9.2 – Tabela solução do método de Ralston

i	x_i	k_1	k_2	$y(x_{i+1})$
0	0,0	0,0	0,15	2,02
1	0,2	0,18	0,303	2,0724
2	0,4	0,3276	0,4285	2,1514
3	0,6	0,4486	0,5313	2,2521
4	0,8	0,5478	0,6157	2,3707

Tabela 9.3 – Tabela solução do método do ponto médio

i	x_i	k_1	k_2	$y(x_{i+1})$
0	0,0	0,0	0,1	2,02
1	0,2	0,18	0,262	2,0724
2	0,4	0,3276	0,3948	2,1514
3	0,6	0,4486	0,5037	2,2521
4	0,8	0,5478	0,5931	2,3707



MÉTODOS DE RUNGE-KUTTA DE TERCEIRA ORDEM - RK3

O método de Runge-Kutta de Terceira Ordem é dado por:

$$y(t_{i+1}) = y(t_i) + h(\alpha_1 k_1 + \alpha_2 k_2 + \alpha_3 k_3) \quad (9.14)$$

com

$$\begin{cases} \alpha_1 + \alpha_2 + \alpha_3 = 1 \\ k_1 = f(t_i, y(t_i)) \\ k_2 = f(t_i + p_1 h, y(t_i) + h(q_{11} k_1)) \quad p_1 = q_{11} \\ k_3 = f(t_i + p_2 h, y(t_i) + h(q_{21} k_1 + q_{22} k_2)) \quad p_2 = q_{21} = q_{22} \end{cases}$$

seguindo o mesmo raciocínio da dedução do método RK2, podemos obter a seguinte expressão:

$$y(t_{i+1}) = y(t_1) + \frac{h}{9} (2k_1 + 3k_2 + 4k_3) \quad (9.15)$$

com

$$\begin{cases} k_1 = f(t_i, y(t_i)) \\ k_2 = f\left(t_i + \frac{h}{2}, y(t_i) + \frac{h}{2}k_1\right) \\ k_3 = f\left(t_i + \frac{3}{4}h, y(t_i) + \frac{3}{4}hk_2\right) \end{cases}$$



que correspondem ao método de **Runge - Kutta de ordem 3** ou **RK3**. e

$$y(t_{i+1}) = y(t_i) + \frac{h}{6} (k_1 + 2k_2 + 2k_3 + k_4) \quad (9.16)$$

com

$$\begin{cases} k_1 = f(t_i, y(t_i)) \\ k_2 = f\left(t_i + \frac{h}{2}, y(t_i) + \frac{h}{2}k_1\right) \\ k_3 = f\left(t_i + \frac{h}{2}, y(t_i) + \frac{h}{2}hk_2\right) \\ k_4 = f(t_i + h, y(t_i) + hk_3) \end{cases}$$

pode ser definido como o método de **Runge - Kutta de ordem 4** ou **RK4**.

Exemplo 9.4

Considerando P como uma determinada população, a equação do modelo de Malthus é dada por,

$$\frac{dP(t)}{dt} = \beta P(t), \quad P(0) = P_0$$

onde β representa a taxa de crescimento intrínseco da população. A solução analítica desta equação é,



$$P(t) = P_0 e^{\beta t}.$$

Apesar de equivocado e criticado em alguns conceitos, este modelo serviu de base à evolução de inúmeros modelos populacionais. A seguir, iremos aplicar técnicas estudadas ao longo da disciplina para resolver um problema de ordem prática e aplicado a tal modelo. Seja o seguinte conjunto de dados, que constitui o senso demográfico em uma determinada cidade ao longo dos anos apresentado:

t_i	0,0	1,0	2,0	3,0	4,0
Ano	1980	1991	2000	2007	2010
População	9204	11620	16197	18687	19825

Solução. Primeiramente iremos ajustar uma curva ao modelo: $P(t) = P_0 e^{\beta t}$.

$$P(t) = P_0 e^{\beta t} \rightarrow \ln(P(t)) = \ln(P_0) + \beta t.$$

Considerando o caso geral temos,

$$\begin{bmatrix} 5 & 10 \\ 10 & 30 \end{bmatrix} \begin{bmatrix} \ln(P_0) \\ \beta \end{bmatrix} = \begin{bmatrix} 47,9107 \\ 97,8312 \end{bmatrix} \Rightarrow \begin{bmatrix} \ln(P_0) \\ \beta \end{bmatrix} = \begin{bmatrix} 9,1802 \\ 0,2 \end{bmatrix} \rightarrow \begin{cases} P_0 = 9703,15 \\ \beta = 0,2 \end{cases}$$

Assim, o modelo ajustado corresponde a $P(t) = 9703,19 e^{0,2t}$ e a EDO inicial pode ser re escrita como:

$$\frac{dP(t)}{dt} = 0,2 \cdot P(t), \quad P_0 = 9703,19.$$

Para resolver numericamente, primeiro iremos utilizar o método de Euler. Considerando $t = 0, 1, 2, 3$. Obtemos:



$$P(t_1) = P(t_0) + 0,2.P(t_0) = 9703,19 + 0,2.9703,19 = 11643,83$$

$$P(t_2) = P(t_1) + 0,2.P(t_1) = 11643,83 + 0,2.11643,83 = 13972,59$$

$$P(t_3) = P(t_2) + 0,2.P(t_2) = 13972,59 + 0,2.13972,59 = 16767,11$$

$$P(t_4) = P(t_3) + 0,2.P(t_3) = 16767,11 + 0,2.16767,11 = 20120,53.$$

Em seguida, considerando as mesmas condições e utilizando o método de Euler modificado. Obtemos:

$$t = 1 \begin{cases} k_1 = 0,2.P(t_1) = 2367,58 \\ k_2 = 0,2.(P(t_1) + h.k_1) = 2841,09 \\ P(t_2) = P(t_1) + \frac{h}{2}(k_1 + k_2) = 14442,22 \end{cases}$$

$$t = 2 \begin{cases} k_1 = 0,2.P(t_2) = 2888,44 \\ k_2 = 0,2.(P(t_2) + h.k_1) = 3466,13 \\ P(t_3) = P(t_2) + \frac{h}{2}(k_1 + k_2) = 17619,52 \end{cases}$$

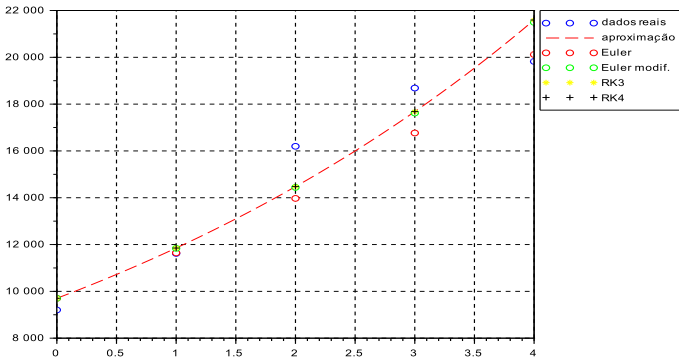
$$t = 3 \begin{cases} k_1 = 0,2.P(t_3) = 3523,90 \\ k_2 = 0,2.(P(t_3) + h.k_1) = 4228,68 \\ P(t_4) = P(t_3) + \frac{h}{2}(k_1 + k_2) = 21495,81 \end{cases}$$

Para finalizar apresentaremos uma tabela de valores obtidos utilizando diferentes métodos de solução numérica para EDO's:

t_i	sol. analítica	Ajuste	Euler	Euler modif.	RK3	RK4
0	9204	9703,19	9703,19	9703,19	9703,19	9703,19
1	11620	11851,50	11643,83	11837,89	11850,83	11851,48
2	16197	14475,45	13972,59	14442,22	14473,81	14475,40
3	18687	17680,36	16767,11	17619,52	17677,35	17680,24
4	19825	21594,85	20120,53	21495,81	21589,94	21594,65



Figura 9.2 – Dados ajustados e suas respectivas soluções numéricas



EDO'S DE ORDEM SUPERIOR

Uma EDO de ordem superior (maior que 1), dada por

$$y^{(n)} = f(t, y, y', y'', \dots, y^{(n-1)})$$

de modo que a EDO original pode ser decomposta como um sistema de n equações diferenciais de primeira ordem. Ou seja,

$$\begin{cases} y_1' = f_1(t, y_1, y_2, \dots, y_n) \\ y_2' = f_2(t, y_1, y_2, \dots, y_n) \\ \vdots \\ y_n' = f_n(t, y_1, y_2, \dots, y_n) \end{cases}$$



em que

$$\begin{cases} y_1 = y \\ y_2 = y' \\ \vdots \\ y_n = y^{(n-1)} \end{cases}$$

com f_1, f_2, \dots, f_n funções dadas e as condições iniciais:

$$y_1(t_0) = y_0, y_2(t_0) = y'_0, \dots, y_n(t_0) = y_0^{(n-1)}$$

que pode ser escrito na forma matricial e facilmente ser resolvido por qualquer método numérico de resolução de EDO's de primeira ordem.

Por exemplo, o método de Euler pode ser escrito vetorialmente da seguinte forma:

$$\mathbf{Y}_{i+1} = \mathbf{Y}_i + h\mathbf{F}(t_i, \mathbf{Y}_i)$$

Seja o sistema modelado pela seguinte EDO,

$$\frac{d^2\theta}{dt^2} + \frac{g}{l}\text{sen}(\theta) = 0$$

que pode ser escrita como,

$$\begin{cases} y'_1 = y_2 \\ y'_2 = -\frac{g}{l}\text{sen}(y_1) \end{cases}$$

considerando $l = 1m$ e as condições iniciais,



$$y_1(0) = \frac{\pi}{2}, y_2(0) = 0.$$

EXERCÍCIOS PROPOSTOS

Exercício 9.1

Considere a EDO de primeira ordem:

$$\frac{dy}{dx} = yx - x^3$$

com $0 \leq x \leq 1, 8$ e $y(0) = 1$. Resolva manualmente utilizando o método de Euler com

Exercício 9.2

Considere o seguinte sistema de duas equações diferenciais:

$$\frac{dx}{dt} = x + y$$

$$\frac{dy}{dt} = y - x$$

com $0 \leq t \leq 2$, $x(0) = 1$ e $y(0) = 3$.



- Resolva manualmente utilizando o método de Euler e o RK4 com $h = 0,5$.
- A solução analítica do sistema é $x(t) = e^t (3\text{sen}(t) + \text{cos}(t))$ e $y(t) = e^t (3\text{cos}(t) - \text{sen}(t))$. Calcule em cada passo o erro absoluto entre a solução exata e a solução numérica.

Exercício 9.3

Se água for drenada de um tanque cilíndrico vertical abrindo-se uma válvula na base, ela escoará rapidamente quando o tanque estiver cheio e mais lentamente conforme ele continuar a ser drenado. Como pode ser mostrado, a taxa pela qual o nível de água decresce é:

$$\frac{dh}{dt} = -k\sqrt{h}$$

em que k é uma constante dependente da forma do orifício, da área da seção transversal do tanque e do orifício de drenagem. A profundidade da água h é medida em metros e o tempo t , em minutos. Se $k = 0,06$, determine quanto tempo levará para que o tanque fique vazio se o nível de fluido for inicialmente 3 m . Use um passo de meio minuto.

Exercício 9.4

Supondo que o arrasto seja proporcional ao quadrado da velocidade, podemos modelar a velocidade de um objeto em queda livre como o pára-quedista com a seguinte equação diferencial:



$$\frac{dv}{dt} = g - \frac{c_d}{m}v^2$$

em que v é a velocidade (m/s), t é o tempo (s), g é a aceleração da gravidade ($\approx 9,81 \text{ m/s}^2$), c_d é um coeficiente de arrasto de segunda ordem (kg/m) e m é a massa (kg). Resolva para determinar a velocidade e a distância percorrida na queda por um objeto de 90 kg com um coeficiente de arrasto de 0,225 kg/m. Se a altura inicial for 1 km, determine quando ele atinge o chão.

Exercício 9.5

Um Boeing de massa $m = 97 \times 10^3 \text{ kg}$ aterrissa a uma velocidade de 93 m/s e liga os seus reversos em $t = 0$. A força F aplicada no avião à medida que ele reduz a sua velocidade é dada por $F = -(5v^2 + 57 \times 10^4)$, onde v é a velocidade do avião. Usando a segunda lei de Newton do movimento e da dinâmica dos fluidos, a relação entre a velocidade e a posição x do avião pode ser escrita como:

$$mv \frac{dv}{dx} = F$$

onde x é a distância medida a partir da localização do jato em $t = 0$. Determine a distância percorrida pelo avião antes que sua velocidade se reduza a 40 m/s usando um dos métodos de resolução numérica de integrais estudado.



Em seguida, resolva a EDO considerando $0 \leq x \leq 400 \text{ m}$ e $\Delta x = 50 \text{ m}$.

Exercício 9.6

As equações

$$\frac{d^2 r}{dt^2} = \frac{1}{r^2} \left(\frac{9}{r} - 2 \right), \quad \frac{d\theta}{dt} = \frac{3}{r^2}.$$

descrevem a órbita newtoniana de uma partícula em um campo gravitacional, após escolhas adequadas de algumas constantes. Se $t = 0$ na posição em que r é mínimo e $r(0) = 3$, $\theta(0) = 0$ e $r' = 0$. Então a órbita gera uma elipse dada por $r = 9/(2 + \cos \theta)$. Use um dos métodos estudado para obter a solução e compare graficamente com a solução analítica.

Exercício 9.7

Escreva as seguintes EDOs de segunda ordem como um sistema de EDOs de primeira ordem:

- $\frac{d^2 y}{dx^2} = -Py + \frac{QL}{2}x - \frac{Q}{2}x^2$, onde P, Q, L são constantes.
- $\frac{d^2 y}{dx^2} = M \left[1 + \left(\frac{dy}{dx} \right)^2 \right]^{\frac{3}{2}}$, onde M é constante.



Exercício 9.8

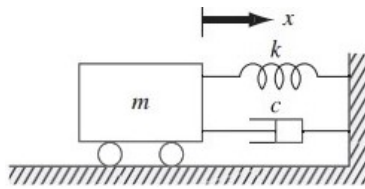
Um indutor e um resistor não-linear de resistência $R = 500 + 250I^2 \Omega$ estão conectados em série com uma fonte de tensão CC e uma chave. A chave está inicialmente aberta, sendo então fechada no tempo $t = 0$. A corrente I no circuito para $t > 0$ é determinada a partir da solução da equação:

$$\frac{dI}{dt} = \frac{V_0}{L} - \frac{R}{L}I$$

Para $V_0 = 1000V$ e $L = 15H$, determine e trace a corrente em função do tempo (em segundos).

Exercício 9.9

O movimento de um sistema massa/mola/amortecedor conforme mostrado na seguinte figura:



é descrito pela seguinte equação diferencial,

$$m \frac{d^2x}{dt^2} + c \frac{dx}{dt} + kx = 0$$



em que x é o deslocamento a partir da posição de equilíbrio (m), t é o tempo (s), $m = 20\text{kg}$ é a massa e c é o coeficiente de amortecimento ($\text{N} \cdot \text{s}/\text{m}$). O coeficiente de amortecimento c será analisado com três valores: 5 (subamortecido), 40 (criticamente amortecido) e 200 (superamortecido). A constante da mola é $k = 20\text{N}/\text{m}$. A velocidade inicial é zero e o deslocamento inicial é $x = 1\text{m}$.

Resolva essa equação a partir de um sistema de equações diferenciais de primeira ordem usando um método numérico no período de tempo $0 \leq t \leq 1, 5\text{s}$.

Exercício 9.10

Encontre a equação diferencial que modela o sistema físico do pêndulo (com e sem atrito) e,

- resolva a EDO computacionalmente através dos métodos de Euler e RK4;
- apresente gráficos da velocidade e posição em relação ao tempo para diferentes valores de massa e do comprimento da haste;
- apresente os gráficos do diagrama de fase para os diferentes cenários simulados;
- apresente conclusões coerentes dos resultados obtidos.



UNIDADE IV



ÁLGEBRA LINEAR

Embora um primeiro curso de álgebra linear seja um pré-requisito para a disciplina de computação numérica, para fins de completude, estabelecemos algumas notações e revisamos rapidamente definições e conceitos básicos sobre espaços vetoriais, transformações lineares, normas e matrizes neste apêndice. Resultados fundamentais em normas vetoriais e matriciais são descritos em alguns detalhes. Esses resultados serão usados com frequência ao longo dos capítulos apresentados anteriormente. Os alunos podem revisar este material sempre que necessário. Afinal, a análise de algoritmos matriciais requer o uso de tais ferramentas. Por exemplo, a qualidade de uma solução em sistemas lineares pode não ser a ideal se a matriz dos coeficientes não for adequada. Para isso, precisamos de um quantificador, uma métrica (por exemplo, distância) no espaço matricial. Normas matriciais, por exemplo, podem ser utilizadas para fornecer essa métrica ou quantificador.

ESPAÇOS VETORIAL

Iremos ao longo deste, estudar os espaços de vetores e suas características. Portanto, iniciaremos por meio da definição de nosso principal objeto de estudo.



Definição 10.1

Um espaço vetorial consiste de um conjunto V com elementos denominados de vetores, satisfazendo os seguintes axiomas:

Dados dois elementos $\mathbf{v}, \mathbf{u} \in V$, e uma função

$$+ : V \times V \rightarrow V, \\ (\mathbf{v}, \mathbf{u}) \rightarrow \mathbf{v} + \mathbf{u}$$

denominada soma.

Satisfazendo as seguintes propriedades,

- *comutativa*: $\mathbf{v} + \mathbf{u} = \mathbf{u} + \mathbf{v}, \forall \mathbf{v}, \mathbf{u} \in V$
- *associativa*: $\mathbf{v} + (\mathbf{u} + \mathbf{w}) = (\mathbf{v} + \mathbf{u}) + \mathbf{w}, \forall \mathbf{v}, \mathbf{u}, \mathbf{w} \in V$
- *elemento neutro*: $\exists! \mathbf{0} \in V$ tal que $\mathbf{v} + \mathbf{0} = \mathbf{v} \forall \mathbf{v} \in V$
- *inverso aditivo*: $\exists! -\mathbf{v} \in V$ tal que $\mathbf{v} + (-\mathbf{v}) = \mathbf{0} \forall \mathbf{v} \in V$

Dados qualquer elemento $\mathbf{v} \in V$, um escalar $\lambda \in \mathbb{R}$ e uma função.

$$\mathbb{R} \times V \rightarrow V, \\ (\lambda, \mathbf{v}) \rightarrow \lambda \cdot \mathbf{v}$$

denominada produto por um escalar. Satisfazendo as seguintes propriedades,

- *associativa*: $\alpha(\lambda \mathbf{v}) = (\alpha\lambda) \mathbf{v}, \forall \alpha, \lambda \in \mathbb{R} \text{ e } \mathbf{v} \in V$
- *distributiva com relação a escalar*: $\alpha(\mathbf{v} + \mathbf{u}) = \alpha\mathbf{v} + \alpha\mathbf{u}, \forall \alpha \in \mathbb{R} \text{ e } \mathbf{v}, \mathbf{u} \in V$
- *distributiva com relação a vetor*: $(\alpha + \lambda) \mathbf{v} = \alpha\mathbf{v} + \lambda \mathbf{v}, \forall \alpha, \lambda \in \mathbb{R} \text{ e } \mathbf{v} \in V$
- *elemento neutro*: $1 \cdot \mathbf{v} = \mathbf{v} \in V$

A relação entre um espaço vetorial V e um corpo F . É usualmente descrita como sendo um espaço vetorial V sobre o corpo F . Se $F = \mathbb{R}$, temos um *espaço vetorial real*.



Exemplo 10.1

Seja F um corpo. O conjunto F^F de todas as funções de F em F é um espaço vetorial sobre F , cujas operações podem ser vistas como:

$$(f + g)(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$$

$$(\alpha f)(\mathbf{x}) = \alpha f(\mathbf{x}).$$

SUBESPAÇOS VETORIAL

Muitas estruturas algébricas possuem subestruturas, e os espaços vetoriais não são exceção.

Definição 10.2

Um subconjunto não vazio $S \subseteq V$ é dito **subespaço** se qualquer par de vetores, \mathbf{v} e $\mathbf{u} \subseteq S \Rightarrow \alpha \mathbf{v} + \beta \mathbf{u} \subseteq S, \forall \alpha, \beta \in F$.

Definição 10.3

Se V é um espaço vetorial sobre um corpo F , então o conjunto $\{0\}$ e V são subespaços do próprio V .

COMBINAÇÃO LINEAR

Definição 10.4

Seja $S \subseteq V$. Uma **combinação linear** de elementos de S é uma expressão da forma,

$$a_1 \mathbf{v}_1 + a_2 \mathbf{v}_2 + \cdots + a_n \mathbf{v}_n.$$



em que $n > 0$, $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n \in S$, e $a_1, a_2, \dots, a_n \in F$. A combinação linear é dita **trivial** se $a_1 = a_2 = \dots = a_n = 0$; caso contrário é dita **não trivial**.

Uma combinação linear pode ser vista como, uma soma finita de vetores ponderados por escalares. O conjunto de todas as combinações lineares de $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ é um subespaço de S , e o menor subespaço de S . É denominado de subespaço gerado por $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$.

Definição 10.5

Um conjunto de vetores $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m$ em V gera o próprio espaço V se cada vetor $\mathbf{v} \in V$ pode se escrito como uma combinação linear de $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m$.

Definição 10.6

Os vetores $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m$ em V são denominados **linearmente dependentes** se existe uma relação não trivial da forma,

$$\sum_{i=1}^m a_i \mathbf{v}_i = a_1 \mathbf{v}_1 + a_2 \mathbf{v}_2 + \dots + a_m \mathbf{v}_m = 0,$$

e nem todo $a_i, i = 0, 1, \dots, m$ é nulo. Por outro lado, se o conjunto de vetores não é linear dependente, denominamos de **linearmente independente**.

Por definição, um conjunto de vetores linearmente independentes, constitui uma **base** no espaço vetorial.

Exemplo 10.2

Um polinômio da forma $a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n$ constitui uma combinação linear de funções polinomiais e $x^0, x^1, x^2, \dots, x^n$, constitui um base no espaço de funções polinomiais.



Um espaço vetorial é dito de *dimensão finita* se possuir base finita. Se todas as bases de um espaço vetorial V , de dimensão finita, contiver o mesmo número de vetores. Este número é chamada de dimensão do espaço vetorial V e denotado por $\dim(V)$.

Exemplo 10.3

Os vetores $e_1 = (1, 0, \dots, 0)^T$, $e_2 = (0, 1, 0, \dots, 0)^T$, \dots , $e_n = (0, 0, \dots, 1)^T$, formam a base *canônica* e geram o espaço vetorial de dimensão finita \mathbb{R}^n . De modo que dado um vetor $\mathbf{v} = (v_1, v_2, \dots, v_n)^T \in \mathbb{R}^n$. Podemos escrevê-lo como uma combinação linear dos vetores da base, ou seja

$$\mathbf{v} = v_1 \mathbf{e}_1 + v_2 \mathbf{e}_2 + \dots + v_n \mathbf{e}_n = v_1(1, 0, \dots, 0)^T + v_2(0, 1, 0, \dots, 0)^T + \dots + v_n(0, 0, \dots, 1)^T$$

R Em geral,

$$\mathbb{R}^n = \underbrace{\mathbb{R} \times \mathbb{R} \times \dots \times \mathbb{R}}_{n \text{ vezes}}$$

TRANSFORMAÇÕES LINEARES

O objetivo nesta sub seção é estabelecer uma relação direta entre matrizes e transformações lineares, que constitui um dos principais assuntos na teoria em Álgebra Linear e conseqüentemente em Computação Numérica.

Definição 10.7

Sejam V e W espaços vetoriais sobre um corpo F . Um mapeamento $T: V \rightarrow W$ é uma **transformação linear** se,



$$T(\alpha \mathbf{v} + \beta \mathbf{w}) = \alpha T(\mathbf{v}) + \beta T(\mathbf{w})$$

para todos escalares $\alpha, \beta \in \mathbb{F}$ e vetores $\mathbf{v}, \mathbf{w} \in V$. O conjunto de todas as transformações lineares de V para W é denotado por $L(V, W)$.

MATRIZES DA TRANSFORMAÇÃO LINEAR

Seja $T : V \rightarrow W$ uma transformação linear, $B = \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$ uma base em V e $D = \{\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_m\}$ uma base em W . O vetor \mathbf{e}_1 é mapeado por T em algum vetor $T(\mathbf{e}_1) \in W$, assim, cada vetor em W , possui uma expansão

$$\begin{aligned} T(\mathbf{e}_1) &= a_{11}\mathbf{f}_1 + a_{12}\mathbf{f}_2 + \dots + a_{1m}\mathbf{f}_m \\ T(\mathbf{e}_2) &= a_{21}\mathbf{f}_1 + a_{22}\mathbf{f}_2 + \dots + a_{2m}\mathbf{f}_m \\ &\vdots \\ T(\mathbf{e}_n) &= a_{n1}\mathbf{f}_1 + a_{n2}\mathbf{f}_2 + \dots + a_{nm}\mathbf{f}_m \end{aligned}$$

onde os coeficientes a_{ij} , $i = 1, 2, \dots, n$; $j = 1, 2, \dots, m$ definem uma matriz $n \times m$.

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \vdots & \dots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nm} \end{pmatrix}$$

denominada de matriz da transformação linear T em relação as bases B e D . Os coeficientes dos vetores $T(\mathbf{e}_1), T(\mathbf{e}_2), \dots, T(\mathbf{e}_n)$ com relação a base D são os vetores coluna da matriz \mathbf{A} .



Exemplo 10.4

Considere a transformação linear $T: \mathbb{R}^3 \rightarrow \mathbb{R}^3$ definida por

$$T(\mathbf{x}, \mathbf{y}, \mathbf{z}) = (\mathbf{x} - 2\mathbf{y}, \mathbf{y}, \mathbf{x} + \mathbf{y} + 4\mathbf{z}).$$

Então, temos que

$$T \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \\ \mathbf{z} \end{pmatrix} = \begin{pmatrix} 1 & -2 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 4 \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \\ \mathbf{z} \end{pmatrix}$$

com matriz da transformação dada por,

$$\mathbf{A} = \begin{pmatrix} 1 & -2 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 4 \end{pmatrix}$$

Exemplo 10.5

Considere a transformação linear $D: P_2 \rightarrow P_2$ o operador derivada, definido no espaço vetorial polinomial de grau menor ou igual a 2. Seja $B = D = (\mathbf{1}, \mathbf{x}, \mathbf{x}^2)$. Então,

$$D(\mathcal{D}(\mathbf{1}))_{\mathcal{D}} = (\mathbf{0})_{\mathcal{D}} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, D(\mathcal{D}(\mathbf{x}))_{\mathcal{D}} = (\mathbf{1})_{\mathcal{D}} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, D(\mathcal{D}(\mathbf{x}^2))_{\mathcal{D}} = (2\mathbf{x})_{\mathcal{D}} = \begin{pmatrix} 0 \\ 2 \\ 0 \end{pmatrix}$$

com matriz da transformação dada por,

$$\mathbf{A} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \end{pmatrix}.$$



Se considerarmos o polinômio $p(x) = 5 + x + 2x^2$, então

$$D(p(x))_{\mathcal{B}} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 5 \\ 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 1 \\ 4 \\ 0 \end{pmatrix}$$

portanto, $D(p(x)) = 1 + 4x$.

R As definições apresentadas tornam $T(V, W)$ um espaço vetorial. Se $V = W$, usaremos simplesmente a notação $T(V)$ e denominamos tais transformações de *Operadores Lineares em V*.

NORMA NO \mathbb{R}^n

Definição 10.8 — Norma.

Seja $V \subseteq \mathbb{R}^n$ um espaço vetorial. Uma função $\|\cdot\| : V \rightarrow \mathbb{R}$ é denominada **norma** se $\forall \mathbf{x}, \mathbf{y} \in V$ e $\lambda \in \mathbb{R}$, temos:

$$\|\mathbf{x}\| \geq 0; \quad \|\mathbf{x}\| = 0 \Leftrightarrow \mathbf{x} = \mathbf{0}$$

$$\|\lambda \mathbf{x}\| = |\lambda| \cdot \|\mathbf{x}\|$$

$$\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|.$$

Um espaço vetorial V , munido com uma norma é geralmente denominado de **espaço vetorial normado**.

Exemplo 10.6

$(\mathbb{R}^n, \|\cdot\|_2)$ é um espaço vetorial normado. Geralmente, denominado de **espaço euclidiano**.



Considere $\mathbf{x} = (x_1 \ x_2 \ \cdots \ x_n)^T \in \mathbb{R}^n$

$$\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} ; p \geq 1 \rightarrow \text{norma } p$$

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i| = |x_1| + |x_2| + \cdots + |x_n| \rightarrow \text{norma um ou da soma}$$

$$\|\mathbf{x}\|_2 = \left(\sum_{i=1}^n |x_i|^2 \right)^{\frac{1}{2}} = \sqrt{x_1^2 + x_2^2 + \cdots + x_n^2} \rightarrow \text{norma Euclidiana}$$

$$\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i| = \max\{|x_1|, |x_2|, \cdots, |x_n|\} \rightarrow \text{norma infinita ou do máximo}$$

Definição 10.9

Dizemos que duas normas $\|\cdot\|_1$ e $\|\cdot\|_2$ são ditas equivalentes se existirem constantes positivas c_1, c_2 , tais que,

$$c_1 \|\mathbf{v}\|_1 \leq \|\mathbf{v}\|_2 \leq c_2 \|\mathbf{v}\|_1, \quad \forall \mathbf{v} \in \mathbb{V}.$$

Num um espaço de dimensão finita, qualquer duas normas são equivalentes.

Exemplo 10.7

PRODUTO INTERNO

Definição 10.10 — Produto interno.

A função $(\cdot \cdot) : \mathbb{V} \times \mathbb{V} \rightarrow \mathbb{R}$ é denominada de **produto interno** se, $\forall \mathbf{x}, \mathbf{y} \in \mathbb{V}$ e $\lambda \in \mathbb{R}$, temos:



$$\langle \mathbf{x}, \mathbf{x} \rangle \geq 0; \quad \langle \mathbf{x}, \mathbf{x} \rangle = 0 \Leftrightarrow \mathbf{x} = \mathbf{0} \quad (10.4)$$

$$\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle \quad (10.5)$$

$$\langle \mathbf{x} + \lambda \mathbf{y}, \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{z} \rangle + \lambda \langle \mathbf{y}, \mathbf{z} \rangle \quad (10.6)$$

$$\langle \mathbf{x}, \lambda \mathbf{y} + \mathbf{z} \rangle = \lambda \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{z} \rangle. \quad (10.7)$$

R se observarmos as equações (A.1) e (A.4) podemos tirar a seguinte relação

$$\|\mathbf{x}\|_2 = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$$

Dados, $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \in \mathbb{R}^n$ e $\mathbf{y} = (y_1, y_2, \dots, y_n)^T \in \mathbb{R}^n$, temos que

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = x_1 y_1 + x_2 y_2 + \dots + x_n y_n = \sum_{k=1}^n x_k y_k$$

se $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \in \mathbb{C}^n$ (espaço n - vetorial complexo) e $\mathbf{y} = (y_1, y_2, \dots, y_n)^T \in \mathbb{C}^n$, temos que

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y}^* = x_1 y_1^* + x_2 y_2^* + \dots + x_n y_n^* = \sum_{k=1}^n x_k y_k^* = \mathbf{y}^H \mathbf{x}$$

onde \mathbf{y}^* representa o complexo conjugado, \mathbf{y}^H o hermitiano do vetor \mathbf{y} ($\mathbf{y}^H = \mathbf{y}^{*T}$) e, no espaço vetorial \mathbb{C}^n o produto interno é geralmente denominado de **produto hermitiano**.

Definição 10.11 — Ângulo entre vetores.

Seja $V \subseteq \mathbb{R}^n$ um espaço vetorial, o ângulo $\theta \in [0, \pi]$ entre dois vetores não nulos $\mathbf{v}, \mathbf{w} \in V$ é definido como,

$$\cos \theta = \frac{\langle \mathbf{v}, \mathbf{w} \rangle}{\|\mathbf{v}\| \|\mathbf{w}\|} = \frac{\mathbf{v}^T \mathbf{w}}{\|\mathbf{v}\| \|\mathbf{w}\|}. \quad (10.8)$$



Dois vetores no \mathbb{R}^n ou no \mathbb{C}^n são ditos ortogonais se,

$$\langle \mathbf{x}, \mathbf{y} \rangle = 0.$$

Também são ditos ortonormais se,

$$\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 = 1.$$

Além disto, um conjunto de vetores $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\} \in V$ é ortonormal se,

$$\langle \mathbf{e}_i, \mathbf{e}_j \rangle = \begin{cases} 1 & \text{se } i = j, \\ 0 & \text{se } i \neq j. \end{cases} \quad (10.9)$$

um vetor no \mathbb{R}^n pode ser escrito como uma combinação linear de uma base ortonormal, ou seja,

$$\mathbf{x} = x_1 \mathbf{e}_1 + x_2 \mathbf{e}_2 + \dots + x_n \mathbf{e}_n = \langle \mathbf{x}, \mathbf{e}_1 \rangle \mathbf{e}_1 + \dots + \langle \mathbf{x}, \mathbf{e}_n \rangle \mathbf{e}_n$$

$$\|\mathbf{x}\|_2^2 = |\langle \mathbf{x}, \mathbf{e}_1 \rangle|^2 + \dots + |\langle \mathbf{x}, \mathbf{e}_n \rangle|^2 \rightarrow \langle \mathbf{x}, \mathbf{e}_i \rangle = x_i, \quad i = 1, \dots, n$$

onde,

$$\begin{cases} \mathbf{e}_1 = (1, 0, \dots, 0, 0)^T \\ \mathbf{e}_2 = (0, 1, 0, \dots, 0)^T \\ \vdots \\ \mathbf{e}_n = (0, 0, \dots, 0, 1)^T \end{cases}$$



PROJEÇÃO ORTOGONAL

Definição 10.12 — Projeção Ortogonal.

Suponha $V \subset \mathbb{R}^n$ de dimensão m . Dado um vetor $x \in \mathbb{R}^n$, $\exists! y \in V$ tal que $x - y$ é ortogonal a todo vetor de V . O vetor y é denominado de **projeção ortogonal** do vetor x sobre o espaço V .

Exemplo 10.8

Seja, E um espaço vetorial de funções polinomiais no intervalo $[0, 1]$. Para cada $n \in \mathbb{N}$, seja E_n o subespaço de funções polinomiais de grau $\leq n$. Se E_n possui dimensão $n + 1$. Então, $\{p_0, p_1, \dots, p_n\}$ representa uma base do E_n e um polinômio arbitrário de grau n pode ser escrito como uma combinação linear

$$p = a_0 + a_1 p_1 + a_2 p_2 + \dots + a_n p_n = a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n.$$

onde, $a_i \in \mathbb{R}$ e $p_i \in E_n$, para $(i = 0, 1, \dots, n)$. Ou seja,

$$\begin{cases} p_0 = 1 \\ p_1 = x \\ p_2 = x^2 \\ \vdots \\ p_n = x^n \end{cases}$$

Num espaço de funções, podemos definir algumas normas, tais como

Definição 10.13

seja E um espaço de funções contínuas no intervalo $[-1; 1]$. Se



$f, g \in E$, então

$$\|f\|_1 = \int_{-1}^1 |f(x)| dx \rightarrow \text{norma } L^1$$

$$\langle f, g \rangle = \int_{-1}^1 f(x)g(x)dx \rightarrow \text{norma } L^2 \text{ em } \mathbb{E} \subseteq \mathbb{R}^n$$

$$\langle f, g \rangle = \int_{-1}^1 f(x)g^*(x)dx \rightarrow \text{norma } L^2 \text{ em } \mathbb{E} \subseteq \mathbb{C}^n$$

$$\|f\|_2 = \langle f, f \rangle^2 = \left(\int_{-1}^1 f(x)^2 dx \right)^{\frac{1}{2}} \rightarrow \text{norma } L^2$$

R nos exercícios para não sobrecarregar a notação, iremos considerar $\|\cdot\|_2 = \|\cdot\|$

Exercício 10.1

Suponha V um espaço vetorial real com produto interno.

- Mostre que $\langle \mathbf{x} + \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle = \|\mathbf{x}\|^2 - \|\mathbf{y}\|^2, \forall \mathbf{x}, \mathbf{y} \in V$

$$\langle \mathbf{x} + \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle = (\mathbf{x} + \mathbf{y})^T (\mathbf{x} - \mathbf{y}) = \mathbf{x}^T \mathbf{x} - \mathbf{x}^T \mathbf{y} + \mathbf{y}^T \mathbf{x} - \mathbf{y}^T \mathbf{y}$$

$$= \mathbf{x}^T \mathbf{x} - \mathbf{y}^T \mathbf{y} = \langle \mathbf{x}, \mathbf{x} \rangle - \langle \mathbf{y}, \mathbf{y} \rangle = \|\mathbf{x}\|^2 - \|\mathbf{y}\|^2.$$

- Mostre que se $\|\mathbf{x} + \mathbf{y}\| = \|\mathbf{x} - \mathbf{y}\|$, para $\mathbf{x}, \mathbf{y} \in V$. Então \mathbf{x} é ortogonal a \mathbf{y} .



$$\|\mathbf{x} + \mathbf{y}\| = \|\mathbf{x} - \mathbf{y}\| \rightarrow \|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x} - \mathbf{y}\|^2$$

$$\rightarrow \langle \mathbf{x} + \mathbf{y}, \mathbf{x} + \mathbf{y} \rangle = \langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle$$

$$\rightarrow \|\mathbf{x}\|^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle + \|\mathbf{y}\|^2 = \|\mathbf{x}\|^2 - 2\langle \mathbf{x}, \mathbf{y} \rangle + \|\mathbf{y}\|^2 \rightarrow 4\langle \mathbf{x}, \mathbf{y} \rangle = 0$$

$$\rightarrow \langle \mathbf{x}, \mathbf{y} \rangle = 0 \rightarrow \mathbf{x} \perp \mathbf{y}.$$

Exercício 10.2

Suponha $\|\mathbf{x}\| = 3$, $\|\mathbf{x} + \mathbf{y}\| = 5$ e $\|\mathbf{x} - \mathbf{y}\| = 7$. Calcule $\|\mathbf{y}\|$

$$\begin{cases} \|\mathbf{x} + \mathbf{y}\|^2 = \langle \mathbf{x} + \mathbf{y}, \mathbf{x} + \mathbf{y} \rangle = \|\mathbf{x}\|^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle + \|\mathbf{x}\|^2 = 25 \\ \|\mathbf{x} - \mathbf{y}\|^2 = \langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle = \|\mathbf{x}\|^2 - 2\langle \mathbf{x}, \mathbf{y} \rangle + \|\mathbf{x}\|^2 = 49 \end{cases} \rightarrow 2(\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2) = 74.$$

$$\|\mathbf{y}\|^2 = 28 \rightarrow \|\mathbf{y}\| = 2\sqrt{7}$$

Exercício 10.3

Considerando o espaço de funções $C([0, 2\pi]) \subset \mathbb{R}$, com produto interno,

$$\langle f(x), g(x) \rangle = \int_0^{2\pi} f(x)g(x)dx.$$

Mostre que o conjunto $S = \{1, \cos(x), \cos(2x), \dots, \cos(nx), \sin(x), \sin(2x), \dots, \sin(mx)\}$, é ortogonal ($\forall m, n \in \mathbb{Z}$).



$$\langle 1, \cos(nx) \rangle = \int_0^{2\pi} \cos(nx) dx = \frac{1}{n} \operatorname{sen}(nx) \Big|_0^{2\pi} = 0$$

$$\langle 1, \operatorname{sen}(mx) \rangle = \int_0^{2\pi} \operatorname{sen}(mx) dx = -\frac{1}{m} \cos(mx) \Big|_0^{2\pi} = 0$$

$$\langle \operatorname{sen}(mx), \cos(nx) \rangle = \int_0^{2\pi} \operatorname{sen}(mx) \cos(nx) dx = 0.$$

Se $m \neq n$, então o produto das funções trigonométricas, pode ser reescrito como uma soma, ou seja,

$$\begin{cases} \operatorname{sen}((m+n)x) = \operatorname{sen}(mx)\cos(nx) + \operatorname{sen}(nx)\cos(mx) \\ \operatorname{sen}((m-n)x) = \operatorname{sen}(mx)\cos(nx) - \operatorname{sen}(nx)\cos(mx) \\ \operatorname{sen}((m+n)x) + \operatorname{sen}((m-n)x) = 2\operatorname{sen}(mx)\cos(nx) \\ \cos((m+n)x) = \cos(mx)\cos(nx) - \operatorname{sen}(mx)\operatorname{sen}(nx) \\ \cos((m-n)x) = \cos(mx)\cos(nx) + \operatorname{sen}(mx)\operatorname{sen}(nx) \\ \cos((m-n)x) + \cos((m+n)x) = 2\operatorname{sen}(mx)\operatorname{sen}(nx) \end{cases} \quad (10.10)$$

$$\langle \operatorname{sen}(mx), \cos(nx) \rangle = \frac{1}{2} \int_0^{2\pi} \operatorname{sen}((m+n)x) + \operatorname{sen}((m-n)x) dx$$

$$= \frac{1}{2} \left[\frac{\cos((m+n)x)}{m+n} + \frac{\cos((m-n)x)}{m-n} \right]_0^{2\pi} = 0.$$

Se $m = n$, simplesmente teremos



$$\langle \text{sen}(mx), \cos(nx) \rangle = \frac{1}{2m} [\text{sen}^2(mx)]_0^{2\pi} = 0$$

$$\langle \text{sen}(mx), \text{sen}(nx) \rangle = \frac{1}{2} \int_0^{2\pi} \cos((m-n)x) + \cos((m+n)x) dx = 0, \quad m \neq n.$$

$$\langle \cos(mx), \cos(nx) \rangle = \frac{1}{2} \int_0^{2\pi} \cos((m+n)x) + \cos((m-n)x) dx = 0, \quad m \neq n.$$

Note que o conjunto S forma uma base **ortogonal** no espaço $C([0, 2\pi])$. Podemos aproveitar e ortonormalizar esta base, para isso,

$$\left\{ \begin{array}{l} \|1\|^2 = \langle 1, 1 \rangle^2 = \left(\int_0^{2\pi} 1^2 dx \right)^2 = 4\pi \\ \|\text{sen}(mx)\|^2 = \langle \text{sen}(mx), \text{sen}(mx) \rangle^2 = \left(\int_0^{2\pi} \text{sen}^2(mx) dx \right)^2 = 2\pi \\ \|\cos(nx)\|^2 = \langle \cos(nx), \cos(nx) \rangle^2 = \left(\int_0^{2\pi} \cos^2(nx) dx \right)^2 = 2\pi \end{array} \right. \quad (10.11)$$

com isso, o conjunto

$$\left\{ \frac{1}{\sqrt{2\pi}}, \frac{1}{\sqrt{\pi}} \cos(x), \dots, \frac{1}{\sqrt{\pi}} \cos(nx), \frac{1}{\sqrt{\pi}} \text{sen}(x), \dots, \frac{1}{\sqrt{\pi}} \text{sen}(mx) \right\}. \quad (10.12)$$

forma uma base ortonormal.

Exercício 10.4

Sejam



$$\begin{cases} S_k(n) = e^{2\pi i \frac{nk}{N}} \\ S_l(n) = e^{2\pi i \frac{nl}{N}} \end{cases} \quad (10.13)$$

os k, l -ésimos termos complexos de uma função de base. Mostrar que tais funções são ortogonais e, a partir das mesmas, crie um conjunto ortonormal.

ortogonalidade:

$$\begin{aligned} \langle S_k, S_l \rangle &= \sum_{n=0}^{N-1} S_k(n) S_k^*(n) = \sum_{n=0}^{N-1} e^{2\pi i \frac{nk}{N}} e^{-2\pi i \frac{nl}{N}} = \sum_{n=0}^{N-1} e^{2\pi i \frac{n(k-l)}{N}} \\ &= \frac{1 - e^{2\pi i (k-l)}}{1 - e^{2\pi i \frac{(k-l)}{N}}} \rightarrow S_k \perp S_l, k \neq l. \end{aligned}$$

Norma,

$$\langle S_k, S_k \rangle = \sum_{n=0}^{N-1} e^{2\pi i \frac{(k-k)}{N}} = N \rightarrow \|S_k\| = \sqrt{N}.$$

Portanto,



$$\tilde{S}_k(n) = \frac{S_k(n)}{\|S_k\|} = \frac{e^{2\pi i \frac{nk}{N}}}{\sqrt{N}}$$

constitui um conjunto **ortonormal**. Pois,

$$\langle \tilde{S}_k, \tilde{S}_l \rangle = \begin{cases} 1 & \text{se } k = l \\ 0 & \text{se } k \neq l \end{cases}$$

Proposição 10.1

Se $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ formam uma base ortonormal num espaço com produto interno V , então um vetor $\mathbf{w} \in V$, pode ser representado com relação a base por,

$$\mathbf{w} = \langle \mathbf{w}, \mathbf{v}_1 \rangle \mathbf{v}_1 + \langle \mathbf{w}, \mathbf{v}_2 \rangle \mathbf{v}_2 + \dots + \langle \mathbf{w}, \mathbf{v}_n \rangle \mathbf{v}_n. \quad ($$

Teorema 10.1

Seja V um espaço vetorial munido de um produto interno sobre o corpo dos complexos e seja $T: V \rightarrow \mathbb{C}$ um operador (ou funcional) linear em V . Então existe um único vetor $\mathbf{y} \in V$ tal que,

$$T(\mathbf{x}) = \langle \mathbf{x}, \mathbf{y} \rangle \quad \forall \mathbf{x} \in V.$$



Em outras palavras, cada funcional linear é dado por um produto interno.

MATRIZES

Vetores e matrizes são úteis para representar dados numéricos multivariados, e ocorrem naturalmente ao trabalhar com equações lineares ou ao expressar relações lineares entre objetos. Inúmeros algoritmos numéricos envolvem matrizes e operações vetoriais. A grosso modo, uma matriz pode ser vista como um conjunto de vetores linhas ou colunas, onde o número de linhas e colunas determina a dimensão (ou forma) da matriz. Por exemplo, a matriz \mathbf{A} de dimensão $m \times n$ pode ser escrita como,

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ \vdots & \ddots & \cdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}$$

Geralmente utilizamos a notação $\mathbf{A}[i, j]$ para nos referir ao elemento da i -ésima linha e j -ésima coluna da matriz \mathbf{A} ; ou seja, $(\mathbf{A})[i, j] = a_{i, j}$. Se $m = n$ a matriz é dita **quadrada**

Definição 10.14

Seja $\mathbf{A}_{m \times n}$ uma matriz de dimensão $m \times n$. A matriz $\mathbf{B}_{n \times m}$ de dimensão $n \times m$, tal que $b_{ji} = a_{ij}$ é denominada **transposta** da matriz \mathbf{A} . Ou seja, $\mathbf{B} = \mathbf{A}^T$.

- Exemplo 10.9



$$\mathbf{A} = \begin{pmatrix} 2 & 5 & 0 \\ 1 & 3 & 4 \end{pmatrix} \rightarrow \mathbf{B} = \begin{pmatrix} 2 & 1 \\ 5 & 3 \\ 0 & 4 \end{pmatrix} = \mathbf{A}^T$$

Valendo as seguintes propriedades:

- $\mathbf{A}^{TT} = \mathbf{A}$
- linearidade: $(\alpha\mathbf{A} + \gamma\mathbf{B})^T = \alpha\mathbf{A}^T + \gamma\mathbf{B}^T$
- produto: $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$.

ALGUMAS MATRIZES ESPECIAIS

MATRIZES DIAGONAL E TRIANGULARES

Definição 10.15

Uma matriz $\mathbf{A}_{m \times m}$ é dita **diagonal** se $a_{ij} = 0$ para $i \neq j$.

$$\begin{pmatrix} a_{11} & & 0 \\ & \ddots & \\ 0 & & a_{nn} \end{pmatrix}$$

É dita **triangular superior** se $a_{ij} = 0$ para $i > j$. Se $a_{ij} = 0$ para $i < j$ a matriz é triangular inferior.

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ 0 & a_{22} & \cdots & a_{2n} \\ 0 & 0 & \ddots & \vdots \\ 0 & 0 & \cdots & a_{nn} \end{pmatrix}, \begin{pmatrix} a_{11} & 0 & \cdots & 0 \\ a_{21} & a_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix}$$



DETERMINANTE E MATRIZ INVERSA

Definição 10.16

O determinante constitui uma função $\det : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$, que associa a uma matriz quadrada um escalar. Para simplificar os conceitos começaremos exemplificando o determinante de uma matriz elementar 2×2 .

$$\mathbf{A} = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \Rightarrow \det(\mathbf{A}) = ad - bc \quad (10.15)$$

Exemplo 10.10

$$\mathbf{A} = \begin{pmatrix} 2 & -3 \\ 4 & 2 \end{pmatrix} \rightarrow \det(\mathbf{A}) = 2 \cdot 2 - (-3) \cdot 4 = 16.$$

Definição 10.17

Uma matriz quadrada \mathbf{A} é dita invertível se existir uma matriz quadrada \mathbf{B} tal que

$$\mathbf{AB} = \mathbf{BA} = \mathbf{I}$$

Geralmente a matriz inversa de \mathbf{A} é denotada por \mathbf{A}^{-1} . Uma matriz invertível é frequentemente denotada de **não singular**.



Teorema 10.2

Para uma matriz \mathbf{A} quadrada, são equivalentes:

- \mathbf{A} é não singular;
- $\det(\mathbf{A}) \neq 0$;
- posto de \mathbf{A} é completo;
- $N(\mathbf{A}) = \{\mathbf{0}\}$;
- \mathbf{A}^{-1} existe;
- os vetores linhas e colunas de \mathbf{A} são linearmente independentes;
- os autovalores de \mathbf{A} são não nulos;
- $\forall \mathbf{x}, \mathbf{Ax} = \mathbf{0} \rightarrow \mathbf{x} = \mathbf{0}$;
- o sistema $\mathbf{Ax} = \mathbf{b}$ possui solução única.

Para matrizes $n \times n$ (quadradas), iremos abordar outros conceitos, tais como menores e co fatores.

Definição 10.18

Seja $\mathbf{A}_{n \times n}$ o **menor** M_{ij} do elemento a_{ij} é o determinante da matriz obtida, eliminando-se a i ésima linha a j ésima coluna de \mathbf{A} . O **co fator** C_{ij} é dado por:

$$C_{ij} = (-1)^{i+j} M_{ij}.$$

de modo que,



$$\mathbf{C} = \begin{pmatrix} c_{11} & c_{12} & \cdots & c_{1n} \\ c_{21} & c_{22} & \cdots & c_{2n} \\ \vdots & & & \\ c_{n1} & c_{n2} & \cdots & c_{nn} \end{pmatrix}$$

corresponde a matriz de **co fatores** de \mathbf{A} . A transposta dessa matriz é denominada a **matriz adjunta** da matriz \mathbf{A} e é denotada por $adj(\mathbf{A})$. O determinante pode ser calculado utilizando,

$$\det(\mathbf{A}) = \sum_{j=1}^n a_{ij}c_{ij} = a_{11}c_{11} + a_{12}c_{12} + \cdots + a_{1n}c_{1n}$$

Teorema 10.3

Se \mathbf{A} é uma matriz invertível, então

$$\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} adj(\mathbf{A})$$

Exemplo 10.11

$$\mathbf{A} = \begin{pmatrix} 0 & 2 & 1 \\ 3 & 1 & 2 \\ 5 & 1 & 0 \end{pmatrix}$$



$$M_{11} = \begin{vmatrix} 1 & 2 \\ 1 & 0 \end{vmatrix} = -2; M_{12} = \begin{vmatrix} 3 & 2 \\ 5 & 0 \end{vmatrix} = -10; M_{13} = \begin{vmatrix} 3 & 1 \\ 5 & 1 \end{vmatrix} = -2$$

$$M_{21} = \begin{vmatrix} 2 & 1 \\ 1 & 0 \end{vmatrix} = -1; M_{22} = \begin{vmatrix} 0 & 1 \\ 5 & 0 \end{vmatrix} = -5; M_{23} = \begin{vmatrix} 0 & 2 \\ 5 & 1 \end{vmatrix} = -10$$

$$M_{31} = \begin{vmatrix} 2 & 1 \\ 1 & 2 \end{vmatrix} = 3; M_{32} = \begin{vmatrix} 0 & 1 \\ 3 & 2 \end{vmatrix} = -3; M_{33} = \begin{vmatrix} 0 & 2 \\ 3 & 1 \end{vmatrix} = -6$$

$$C_{11} = (-1)^{1+1}M_{11} = -2; C_{12} = (-1)^{1+2}M_{12} = 10; C_{13} = (-1)^{1+3}M_{13} = -2;$$

$$C_{21} = (-1)^{2+1}M_{21} = 1; C_{22} = (-1)^{2+2}M_{22} = -5; C_{23} = (-1)^{2+3}M_{23} = 10;$$

$$C_{31} = (-1)^{3+1}M_{31} = 3; C_{32} = (-1)^{3+2}M_{32} = 3; C_{33} = (-1)^{3+3}M_{33} = -6.$$

Assim,

$$\det(\mathbf{A}) = 0.C_{11} + 2.C_{12} + 1.C_{13} = 18;$$

ou

$$\det(\mathbf{A}) = 3.C_{21} + 1.C_{22} + 2.C_{23} = 18;$$

ou

$$\det(\mathbf{A}) = 5.C_{31} + 1.C_{32} + 0.C_{33} = 18.$$

De modo que,

$$\mathbf{A}^{-1} = \frac{1}{18} \begin{pmatrix} -2 & 1 & 3 \\ 10 & -5 & 3 \\ -2 & 10 & -6 \end{pmatrix}$$



Exemplo 10.12

Considerando o caso 2×2 . Vamos procurar a inversa de

$$\mathbf{A} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

sabendo que,

$$M_{11} = d; M_{12} = c; M_{21} = b \text{ e } M_{22} = a \rightarrow C_{11} = d; C_{12} = -c; C_{21} = -b \text{ e } C_{22} = a.$$

Portanto,

$$\mathbf{C} = \begin{pmatrix} d & -c \\ -b & a \end{pmatrix}$$

sabendo que $\text{adj}(\mathbf{A}) = \mathbf{C}^T$, temos que,

$$\text{adj}(\mathbf{A}) = \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}.$$

Se $\det(\mathbf{A}) \neq 0$, então implica que a matriz é não singular. Assim,

$$\mathbf{A}^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$



Exemplo 10.13

$$\mathbf{A} = \begin{pmatrix} 2 & -3 \\ 4 & 2 \end{pmatrix} \rightarrow \det(\mathbf{A}) = 16$$

$$\mathbf{A}^{-1} = \frac{1}{16} \begin{pmatrix} 2 & 3 \\ -4 & 2 \end{pmatrix}.$$

Ou seja,

$$\mathbf{A}\mathbf{A}^{-1} = \begin{pmatrix} 2 & -3 \\ 4 & 2 \end{pmatrix} \begin{pmatrix} \frac{1}{8} & \frac{3}{16} \\ -\frac{1}{4} & \frac{1}{8} \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = \mathbf{I}$$

$$\mathbf{A}^{-1}\mathbf{A} = \begin{pmatrix} \frac{1}{8} & \frac{3}{16} \\ -\frac{1}{4} & \frac{1}{8} \end{pmatrix} \begin{pmatrix} 2 & -3 \\ 4 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = \mathbf{I}.$$

Exemplo 10.14

Seja

$$\mathbf{A} = \begin{pmatrix} 2 & 1 & 2 \\ 3 & 2 & 2 \\ 1 & 2 & 3 \end{pmatrix}$$

Calcule $\text{adj}(\mathbf{A})$ e \mathbf{A}^{-1} .



$$\text{adj}(\mathbf{A}) = \begin{pmatrix} \begin{vmatrix} 2 & 2 \\ 2 & 3 \end{vmatrix} & -\begin{vmatrix} 3 & 2 \\ 1 & 3 \end{vmatrix} & \begin{vmatrix} 3 & 2 \\ 1 & 2 \end{vmatrix} \\ -\begin{vmatrix} 1 & 2 \\ 2 & 3 \end{vmatrix} & \begin{vmatrix} 2 & 2 \\ 1 & 3 \end{vmatrix} & -\begin{vmatrix} 2 & 1 \\ 1 & 2 \end{vmatrix} \\ \begin{vmatrix} 1 & 2 \\ 2 & 2 \end{vmatrix} & -\begin{vmatrix} 2 & 2 \\ 3 & 2 \end{vmatrix} & \begin{vmatrix} 2 & 1 \\ 3 & 2 \end{vmatrix} \end{pmatrix}^T = \begin{pmatrix} 2 & 1 & -2 \\ -7 & 4 & 2 \\ 4 & -3 & 1 \end{pmatrix}$$

$$\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \text{adj}(\mathbf{A}) = \frac{1}{5} \begin{pmatrix} 2 & 1 & -2 \\ -7 & 4 & 2 \\ 4 & -3 & 1 \end{pmatrix}$$

NORMA MATRICIAL

Definição 10.19

Seja $V \subseteq \mathbb{R}^{m \times n}$ um espaço isomorfo ao \mathbb{R}^{mn} . Uma função $\|\cdot\| : V \rightarrow \mathbb{R}$ é denominada **norma** se $\forall \mathbf{A}, \mathbf{B} \in V$ e $\alpha \in \mathbb{R}$, temos:

$$\|\mathbf{A}\| \geq 0, \forall \mathbf{A} \in V; \|\mathbf{A}\| = 0 \Leftrightarrow \mathbf{A} = \mathbf{0} \quad (10.16)$$

$$\|\alpha \mathbf{A}\| = |\alpha| \|\mathbf{A}\|, \forall \alpha \in \mathbb{R}, \forall \mathbf{A} \in V \quad (10.17)$$

$$\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|, \forall \mathbf{A}, \mathbf{B} \in V. \quad (10.18)$$

Considere $\mathbf{A} = (a_{ij})_{i=1,2,\dots,m; j=1,2,\dots,n}$



$$\|\mathbf{A}\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}| \rightarrow \textit{norma coluna}$$

$$\|\mathbf{A}\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}| \rightarrow \textit{norma linha}$$

$$\|\mathbf{A}\|_2 = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} \rightarrow \textit{norma de Frobenius}$$

que constituem as normas que serão utilizadas ao longo do texto para matrizes.

Definição 10.20

Produto interno entre matrizes.

Sejam $\mathbf{A} \subseteq \mathbb{R}^{n \times m}$ e $\mathbf{B} \subseteq \mathbb{R}^{n \times m}$ matrizes, o **produto interno matricial** é um mapeamento $T: \mathbb{R}^{n \times m} \times \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$ que satisfaz,

$$\langle \mathbf{A}, \mathbf{B} \rangle = \textit{tr}(\mathbf{A}^T \mathbf{B}). \quad (10.19)$$

Uma outra importante propriedade da norma de Frobenius é obtida a partir da definição apresentada do produto interno. Ou seja,

$$\|\mathbf{A}\|_F = \sqrt{\textit{tr}(\mathbf{A}^T \mathbf{A})} = \sqrt{\langle \mathbf{A}, \mathbf{A} \rangle}$$

AUTOVALORES E AUTOVETORES

Seja V um espaço vetorial, e considerando $A: V \rightarrow V$ um operador linear de V , isto é, uma aplicação linear de V em V . Um elemento $\mathbf{v} \in V$ é denominado **autovetor** de \mathbf{A} , se existe um escalar λ tal que seja,

$$\mathbf{A}\mathbf{v} = \lambda \mathbf{v}$$



encontrar os autovalores da matrix \mathbf{A} é resolver a seguinte equação (característica):

$$\det(\mathbf{A} - \lambda \mathbf{I}) = 0$$

e os respectivos autovetores, satisfazem a equação,

$$(\mathbf{A} - \lambda \mathbf{I})\mathbf{v} = \mathbf{0}$$

Teorema 10.4

Se \mathbf{A} é uma matriz **triangular superior, inferior ou diagonal**, então os autovalores de \mathbf{A} são as entradas na diagonal principal e, o determinante pode ser obtido por meio do produto dos respectivos autovalores.

Exemplo 10.15

$$\mathbf{A} = \begin{pmatrix} 2 & 0 & 0 \\ 1 & 3 & 0 \\ 4 & 1 & 5 \end{pmatrix}$$

os autovalores de \mathbf{A} são $\lambda_1 = 2$, $\lambda_2 = 3$ e $\lambda_3 = 5$, respectivamente e $\det(\mathbf{A}) = \lambda_1 \lambda_2 \lambda_3 = 2 \cdot 3 \cdot 5 = 30$.

Se \mathbf{A} é uma matriz 2×2 com entradas reais, então os autovalores de \mathbf{A} podem ser obtidos por meio de,

$$\lambda^2 - \text{tr}(\mathbf{A})\lambda + \det(\mathbf{A}) = 0$$

onde $\text{tr}(\mathbf{A})$ representa o **traço** (soma das entradas na diagonal principal) e $\det(\mathbf{A})$ o **determinante** da matriz.



Exemplo 10.16

$$\mathbf{A} = \begin{pmatrix} 2 & 1 \\ 4 & 3 \end{pmatrix} \rightarrow \lambda^2 - 5\lambda + 2 = 0$$



CÁLCULO DIFERENCIAL

Nas próximas seções apresentaremos alguns conceitos básicos e estabeleceremos algumas propriedades dos espaços \mathbb{R}^n . Tais resultados serão utilizados com frequência ao longo das notas de aula.

Definição 11.1

Sejam $\mathbf{x} \in \mathbb{R}^n$ e $r > 0$. Dizemos que o conjunto

$$B(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n \mid \|\mathbf{x} - \mathbf{y}\| < r\}$$

é denominado **bola aberta** centrada em \mathbf{x} de raio r . E, o conjunto,

$$B(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n \mid \|\mathbf{x} - \mathbf{y}\| \leq r\}$$

é a **bola fechada** de centro \mathbf{x} e raio r . Portanto, um subconjunto $U \subseteq \mathbb{R}^n$ é dito **aberto** em \mathbb{R}^n se, todo ponto em U é centro de alguma bola aberta inteiramente contida em U .

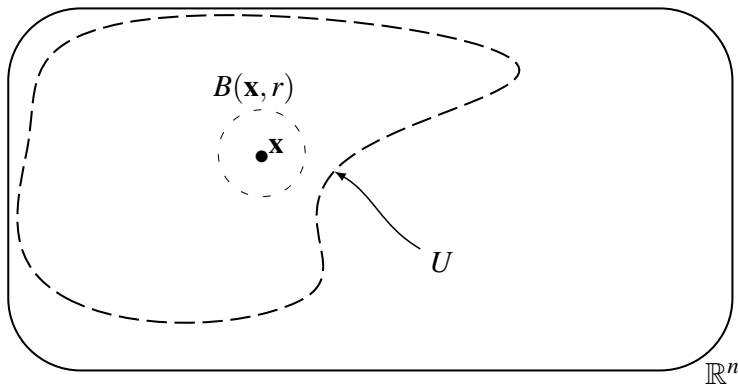
SEQUÊNCIAS CONVERGENTES

Em alguns problemas não é possível a obtenção de uma solução exata. Por isso, na computação numérica trabalhamos geralmente com o conceito de aproximação, e como tal, está sujeito a presença de erros. Portanto, é imprescindível no desenvolvimento dos algoritmos, levar em conta o conceito de convergência.



Uma vez que foi abordado o conceito de norma anteriormente, se faz necessário abordar o conceito de seqüências no \mathbb{R}^n e, sempre que possível estender as mesmas os principais conceitos e resultados das seqüências em \mathbb{R} .

Figura 11.1 – Espaço métrico \mathbb{R}^n com um conjunto aberto $U \subset \mathbb{R}^n$ e uma bola aberta centrada em $\mathbf{x} \in U$ de raio r



Definição 11.2

Seja $S \subseteq \mathbb{R}^n$ um conjunto qualquer, uma seqüência em S é uma função

$$X : \mathbb{N} \rightarrow \mathbb{R}^n \\ n \rightarrow x_n$$

que associa a cada número natural $n \in \mathbb{N}$ um único elemento em \mathbb{R}^n .

Definição 11.3

Seja $V \subseteq \mathbb{R}^n$ um espaço vetorial com uma norma $\|\cdot\|$. Uma seqüência $x_k \in V$ é dita convergente para $\mathbf{x} \in V$ se,

$$\forall \varepsilon > 0, \exists k_0 \in \mathbb{N}, \text{ tal que } k \geq k_0 \rightarrow \|\mathbf{x}_k - \mathbf{x}\| < \varepsilon.$$



Ou seja,

$$\lim_{k \rightarrow \infty} \|\mathbf{x}_k - \mathbf{x}\| = 0.$$

caso contrário, a sequência é dita divergente.

No caso matricial,

Definição 11.4

Uma sequência $(\mathbf{A}_k)_{k \in \mathbb{N}}$ no $\mathbb{R}^{m \times n}$ converge, se existe uma matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ tal que,

$$\lim_{k \rightarrow \infty} \|\mathbf{A}_k - \mathbf{A}\| = 0$$

Definição 11.5

Sejam $U \subseteq \mathbb{R}^n$ e $f: U \rightarrow W$ um mapeamento de U num espaço vetorial normado. Seja $\mathbf{v} \in U$. Dizemos que f é **contínua** em \mathbf{v} se,

$$\lim_{\mathbf{x} \rightarrow \mathbf{v}} f(\mathbf{x}) = f(\mathbf{v}).$$

Ou seja, dado $\varepsilon > 0, \exists \delta > 0 \Rightarrow \|\mathbf{x} - \mathbf{v}\| < \delta \rightarrow \|f(\mathbf{x}) - f(\mathbf{v})\| < \varepsilon$.

Uma função $f: V \rightarrow \mathbb{R}$ é dita *contínua* em $\mathbf{v} \in V$ se para uma sequência $\mathbf{v}_k \rightarrow \mathbf{v}$, temos que $f(\mathbf{v}_k) \rightarrow f(\mathbf{v})$ quando $k \rightarrow \infty$. A função f é dita *contínua em V* se for *contínua para todo $\mathbf{v} \in V$* .

DERIVADAS EM ESPAÇOS VETORIAIS

Sejam U, W espaços vetoriais normados e $T: U \rightarrow W$ uma aplicação linear. Portanto, as seguintes condições em T são equivalentes.

- (1) T é contínua.
- (2) Então existe $\alpha > 0$ tal que $\forall \mathbf{v} \in U$ temos



$$\|T(\mathbf{v})\| \leq \alpha \|\mathbf{v}\|.$$

De fato, se assumirmos (2), então podemos encontrar para todo $\mathbf{v}, \mathbf{u} \in U$

$$\|T(\mathbf{v}) - T(\mathbf{u})\| = \|T(\mathbf{v} - \mathbf{u})\| \leq \alpha \|\mathbf{v} - \mathbf{u}\|.$$

Concluindo que T é sempre uniformemente contínua.

Definição 11.6

A derivada como uma aplicação linear.

sejam $U \subseteq \mathbb{R}^n$ aberto e $\mathbf{v} \in U$. Se existir uma aplicação linear $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ tal que,

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{\|f(\mathbf{v} + \mathbf{h}) - f(\mathbf{v}) - T\mathbf{h}\|}{\|\mathbf{h}\|} = 0, \quad (11.1)$$

dizemos que f é **diferenciável** em \mathbf{v} e denotamos por

$$f'(\mathbf{v}) = T.$$

Se f é diferenciável para todo $\mathbf{v} \in U$, dizemos que f é *diferenciável* em U .

Na equação (11.1), admite-se que $\mathbf{h} \in \mathbb{R}^n$, se $\|\mathbf{h}\| \rightarrow 0$, então $\mathbf{v} + \mathbf{h} \in U$, pois U é aberto. Assim, $f(\mathbf{v} + \mathbf{h}) \in \mathbb{R}^m$ e, como $T \in L(\mathbb{R}^n, \mathbb{R}^m)$, $T\mathbf{h} \in \mathbb{R}^m$. Logo,

$$f(\mathbf{v} + \mathbf{h}) - f(\mathbf{v}) - T\mathbf{h} \in \mathbb{R}^m.$$

Lembrando que, a norma do numerador na equação (11.1) é a de \mathbb{R}^m e no denominador temos a norma de \mathbb{R}^n para \mathbf{h} .



Seja o mapeamento linear,

$$T : (U_1 \times U_2) \subseteq \mathbb{R}^2 \rightarrow (W_1 \times W_2) \subseteq \mathbb{R}^2$$

onde U_1, U_2, W_1 e W_2 são abertos do \mathbb{R}^2 e cada $T_{ij} : U_j \rightarrow W_i$ é, por si só, um mapeamento linear. Seja $\mathbf{v} = (v_1, v_2)^T \in U_1 \times U_2$, que pode ser escrito da forma,

$$\mathbf{v} = v_1 \mathbf{e}_1 + v_2 \mathbf{e}_2,$$

onde $\{\mathbf{e}_1, \mathbf{e}_2\}$ constitui uma base canônica no \mathbb{R}^2 . Portanto,

$$T(\mathbf{v}) = T(v_1, v_2) = T[T_1(v_1, v_2), T_2(v_1, v_2)] \in W_1 \times W_2.$$

Tal que,

$$T(v_1, v_2) = [T_1((v_1, 0) + T_1(0, v_2)), T_2((v_1, 0) + T_2(0, v_2))]$$

onde o mapeamento $v_1 \rightarrow T_1(v_1, 0)$ é linear de $U_1 \subseteq \mathbb{R} \rightarrow W_1 \subseteq \mathbb{R}$, que podemos simplesmente denominar de $T_{11}(v_1)$ ou a_{11} . Similarmente, temos:

$$\begin{cases} T_2(v_1, 0) = T_{21}(v_1) = a_{21} \\ T_2(0, v_2) = T_{12}(v_2) = a_{12} \\ T_2(0, v_2) = T_{22}(v_2) = a_{22}. \end{cases}$$

Então, podemos representar T por meio da seguinte matriz,

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$$



como já discutido na seção sobre transformações lineares, podemos ver que $T(v_1, v_2)$ é dada por meio da multiplicação da matriz acima, com o vetor $\mathbf{v} = (v_1, v_2)^T$. Finalmente podemos observar que se todos $T_{i,j}$ são contínuos.

Se aplicarmos este caso para as derivadas parciais, podemos ter que, dados U um aberto no \mathbb{R}^2 e $f: U \rightarrow F_1 \times F_2$ um mapeamento de classe C^p . Seja $f = (f_1, f_2)$ representada por seus mapeamentos coordenados,

$$f_1: U \rightarrow F_1 \text{ e } f_2: U \rightarrow F_2.$$

Então, para algum $\mathbf{x} = (x_1, x_2)^T \in U$, o mapeamento linear $Df(\mathbf{x})$ é representado por meio da matriz,

$$\begin{pmatrix} D_1 f_1(\mathbf{x}) & D_2 f_1(\mathbf{x}) \\ D_1 f_2(\mathbf{x}) & D_2 f_2(\mathbf{x}) \end{pmatrix}$$

Ou,

$$\begin{pmatrix} \partial f_1 / \partial x_1 & \partial f_1 / \partial x_2 \\ \partial f_2 / \partial x_1 & \partial f_2 / \partial x_2 \end{pmatrix}$$

que corresponde a matriz **Jacobiana** da f em \mathbf{x} e,

$$\begin{pmatrix} \partial f_1 / \partial x_1 & \partial f_1 / \partial x_2 \end{pmatrix} = \nabla f_1(\mathbf{x})$$

$$\begin{pmatrix} \partial f_2 / \partial x_1 & \partial f_2 / \partial x_2 \end{pmatrix} = \nabla f_2(\mathbf{x})$$

correspondem aos vetores gradiente de f_1 e f_2 , respectivamente.



BIBLIOGRAFIA

LIVROS

CHAPRA, Steven C; CANALE, Raymond P. *Métodos numéricos para engenharia*. São Paulo: McGraw-Hill, c2008. xxi, 809 p. ISBN: 9788586804878.

FRANCO, Neide Maria Bertoldi. *Cálculo numérico*. São Paulo: Pearson Prentice Hall, 2006. 505 p. ISBN: 9788576050872.

ARENALES, Selma; *Cálculo Numérico: aprendizagem com apoio de software* 2. ed. São Paulo: Cengage Learning, 2015. ISBN: 9788522112876

EPPERSON, James F., author. *An introduction to numerical methods and analysis* / James F. Epperson, Mathematical Reviews. — Second edition.

VIEIRA, Vandenberg Lopes. *Um curso básico em teoria dos números*. São Paulo: Livraria da Física, 2015.

HALMOS, Paul R. *Finite Dimensional Vector Spaces*. ISBN: 13: 978-0-486-81486-5 9788586804878.

LANG, S. *Undergraduate Analysis*. ISBN: 94841-7.

LANG, S. *Linear Algebra (Undergraduate texts in mathematics)*.

LIMA, Elon L. *Análise Real vol I - Funções de uma variável*. 11. ed. Rio de Janeiro: IMPA, 2011.

LIMA, Elon L. *Álgebra Linear - Coleção Matemática Universitária*, 8. ed. Rio de Janeiro: IMPA, 2009.

AXLER, S. *Linear Algebra - Done Right*.



ANTON, H. e Rorres, C. *Elementary linear algebra with applications.*

SHILOV, Georgi E. *Linear Algebra.*

Atkinson, K.; Han, W. *Theoretical Numerical Analysis - A functional Analysis Framework,*

ARTIGOS

Berrut, Jean-Paul; Trefethen Lloyd N. - *Barycentric Lagrange Interpolation* - SIAM Review, vol. 46, Number 3, pp. 501 - 517

Viana, Geraldo V. R. - Padrão IEEE 754 para Aritmética Binária de Ponto Flutuante.





Composto na
CAULE DE PAPIRO GRÁFICA E EDITORA
Rua Serra do Mel, 7989, Cidade Satélite
Pitimbu | Natal/RN | (84) 3218 4626
cauledepapiro.com.br





editora
CAULE DE PAPIRO®

ISBN 978-65-5477-062-0



9 786554 770620 >