

Matemática Numérica Avançada

Pedro H A Konzen

21 de fevereiro de 2022

Licença

Este trabalho está licenciado sob a Licença Atribuição-CompartilhaIgual 4.0 Internacional Creative Commons. Para visualizar uma cópia desta licença, visite http://creativecommons.org/licenses/by-sa/4.0/deed.pt_BR ou mande uma carta para Creative Commons, PO Box 1866, Mountain View, CA 94042, USA.

Prefácio

Nestas notas de aula são abordados métodos numéricos aplicados a problemas de grande porte. Como ferramenta computacional de apoio, exemplos de aplicação de códigos [Python](#), são apresentados, mais especificamente, códigos com suporte das bibliotecas [NumPy](#) e [SciPy](#).

Agradeço a todos e todas que de modo assíduo ou esporádico contribuem com correções, sugestões e críticas. :)

Pedro H A Konzen

Sumário

Capa	i
Licença	ii
Prefácio	iii
Sumário	iv
1 Sistemas Lineares	1
1.1 Matrizes Esparsas	1
1.1.1 Sistemas Tridiagonais	3
1.1.2 Matrizes Banda	6
1.1.3 Esquemas de Armazenamento	11
1.2 Métodos Iterativos	16
1.2.1 GMRES	16
1.2.2 Método do Gradiente Conjugado	20
Referências Bibliográficas	23

Capítulo 1

Sistemas Lineares

[Vídeo] | [Áudio] | [\[Contatar\]](#)

Neste capítulo, apresentam-se métodos numéricos para a resolução de sistemas lineares de grande porte. Salvo explicitado ao contrário, assume-se que os sistemas são quadrados e têm solução única.

1.1 Matrizes Esparsas

[Vídeo] | [Áudio] | [\[Contatar\]](#)

Uma matriz é dita ser **esparsa** quando ela tem apenas poucos elementos não nulos. A ideia é que os elementos não nulos não precisam ser guardados na memória do computador, gerando um grande benefício na redução da demanda de armazenamento de dados. O desafio está no desenvolvimento de estruturas de dados para a alocação eficiente de tais matrizes, i.e. que sejam suficientemente adequadas para os métodos numéricos conhecidos.

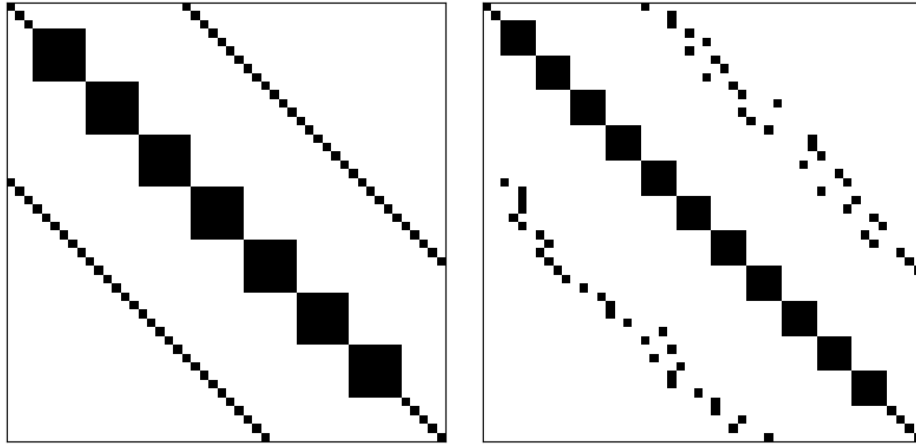


Figura 1.1: Esquerda: exemplo de uma matriz esparsa estruturada. Direita: exemplo de uma matriz esparsa não-estruturada.

Matrizes esparsas podem ser classificadas como **estruturadas** ou **não-estruturadas**. Uma matriz estruturada é aquela em que as entradas não-nulas formam um padrão regular. Por exemplo, estão dispostas em poucas diagonais ou formam blocos (submatrizes densas) ao longo de sua diagonal principal. No caso de não haver um padrão regular das entradas não-nulas, a matriz esparsa é dita ser não-estruturada. Consulte a Figura 1.1 para exemplos.

A **esparsidade** de uma matriz é a porcentagem de elementos nulos que ela tem, i.e. para uma matriz quadrada $n \times n$ tem-se que a esparsidade é

$$\frac{n_{\text{nulos}}}{n^2} \times 100\% \quad (1.1)$$

Por exemplo, a matriz identidade de tamanho $n = 100$ tem esparsidade

$$\frac{100^2 - 100}{100^2} \times 100\% = 99\% \quad (1.2)$$

1.1.1 Sistemas Tridiagonais

Um sistema tridiagonal tem a seguinte forma matricial

$$\begin{bmatrix} b_1 & c_1 & & & 0 \\ a_2 & b_2 & c_2 & & \\ & a_3 & b_3 & \ddots & \\ & & \ddots & \ddots & c_{n-1} \\ 0 & & & a_n & b_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} d_1 \\ d_2 \\ d_3 \\ \vdots \\ d_n \end{bmatrix} \quad (1.3)$$

com, $a_1 = 0$ e $c_n = 0$. Ou seja, é um sistema cuja a matriz dos coeficientes é tridiagonal.

Uma **matriz tridiagonal** é uma matriz esparsa estruturada. Mais especificamente, é chamada de **matriz banda**, em que os elementos não nulos estão apenas em algumas de suas diagonais. Para armazenarmos tal matriz precisamos alocar apenas os seguintes três vetores

$$a = (0, a_2, \dots, a_n) \quad (1.4)$$

$$b = (b_1, b_2, \dots, b_n) \quad (1.5)$$

$$c = (c_1, c_2, \dots, c_{n-1}, 0) \quad (1.6)$$

Ou seja, precisamos armazenar $3n$ pontos flutuantes em vez de n^2 , como seria o caso se a matriz dos coeficientes fosse densa.

Algoritmo de Thomas

O Algoritmo de Thomas¹ é uma forma otimizada do Método de Eliminação Gaussiana aplicada à sistemas tridiagonais. Enquanto este requer $O(n^3)$ operações, esse demanda apenas $O(n)$.

Eliminando os termos abaixo da diagonal em (1.3), obtemos o sistema equivalente

$$\begin{bmatrix} \tilde{b}_1 & c_1 & & & 0 \\ & \tilde{b}_2 & c_2 & & \\ & & \tilde{b}_3 & \ddots & \\ & & \ddots & \ddots & c_{n-1} \\ 0 & & & \tilde{b}_n & \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \tilde{d}_1 \\ \tilde{d}_2 \\ \tilde{d}_3 \\ \vdots \\ \tilde{d}_n \end{bmatrix} \quad (1.7)$$

¹Llewellyn Hilleth Thomas, 1903 - 1992, físico e matemático aplicado britânico. Fonte: [Wikipedia](#).

Este é obtido pela seguinte iteração

$$w := \frac{a_i}{b_{i-1}} \quad (1.8)$$

$$b_i := b_i - wc_{i-1} \quad (1.9)$$

$$d_i := d_i - wd_{i-1} \quad (1.10)$$

onde, o \sim foi esquecido de propósito, indicando a reutilização dos vetores b e d . A solução do sistema é, então, obtida de baixo para cima, i.e.

$$x_n = \frac{d_n}{b_n} \quad (1.11)$$

$$x_i = \frac{d_i - c_i x_{i+1}}{b_i}, \quad (1.12)$$

com $i = n-1, n-2, \dots, 1$.

Listing 1.1: Algoritmo de Thomas

```

1 import numpy as np
2
3 def TDMA(a,b,c,d):
4     n = b.size
5     for i in np.arange(1,n):
6         w = a[i]/b[i-1]
7         b[i] = b[i] - w*c[i-1]
8         d[i] = d[i] - w*d[i-1]
9     x = np.empty(n)
10    x[n-1] = d[n-1]/b[n-1]
11    for i in np.arange(n-2,-1,-1):
12        x[i] = (d[i] - c[i]*x[i+1])/b[i]
13    return x

```

Exercício 1.1.1. Considere o seguinte sistema linear

$$2x_1 - x_2 = 0 \quad (1.13)$$

$$x_{i-1} - 3x_i + 4x_{i+1} = \sin\left(i\frac{\pi}{2(n-1)}\right) \quad (1.14)$$

$$x_{n-1} + x_n = 1 \quad (1.15)$$

a) Compute sua solução usando o Algoritmo de Thomas para $n = 3$.

- b) Compare a solução obtida no item anterior com a gerada pela função `scipy.linalg.solve`.
- c) Compare a solução com a obtida no item anterior com a gerada pela função `scipy.linalg.solve_banded`.
- d) Use o módulo `Python timeit` para comprar a demanda de tempo computacional de cada um dos métodos acima. Compute para $n = 10, 100, 1000, 10000$.

Exercício 1.1.2. Considere que o problema de valor de contorno (PVC)

$$-u'' = \text{sen } \pi x, \quad 0 < x < 1, \quad (1.16)$$

$$u(0) = 0, \quad (1.17)$$

$$u(1) = 0 \quad (1.18)$$

seja simulado com o Método das Diferenças Finitas². Vamos assumir uma discretização espacial uniforme com n nodos e tamanho de malha

$$h = \frac{1}{n-1}. \quad (1.19)$$

Com isso, temos os nodos $x_i = (i-1)h$, $i = 1, 2, \dots, n$. Nos nodos internos, aplicamos a fórmula de diferenças central

$$u''(x_i) \approx \frac{u_{i-1} - 2u_i + u_{i+1}}{h^2}, \quad (1.20)$$

onde, $u_i \approx u(x_i)$. Com isso, a discretização da EDO fornece

$$-\frac{1}{h^2}u_{i-1} + \frac{2}{h^2}u_i - \frac{1}{h^2}u_{i+1} = \text{sen } \pi x_i \quad (1.21)$$

para $i = 2, 3, \dots, n-1$. Das condições de contorno temos $u_1 = u_n = 0$. Logo, o problema discreto lê-se: encontrar $u = (u_1, u_2, \dots, u_n) \in \mathbb{R}^n$ tal que

$$u_1 = 0 \quad (1.22)$$

$$-\frac{1}{h^2}u_{i-1} + \frac{2}{h^2}u_i - \frac{1}{h^2}u_{i+1} = \text{sen } \pi x_i \quad (1.23)$$

$$u_n = 0 \quad (1.24)$$

²Consulte mais em [Notas de Aula: Matemática Numérica](#).

- a) Calcule a solução analítica do PVC.
- b) Use a função `scipy.linalg.solve_banded` para computar a solução do problema discreto associado para diferentes tamanhos de malha $h = 10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}$. Compute o erro da solução discreta em relação à solução analítica.
- c) Compare a demanda de tempo computacional se a função `scipy.linalg.solve` for empregada na computação da solução discreta.

1.1.2 Matrizes Banda

Uma matriz banda é aquela em que os elementos não nulos estão dispostos em apenas algumas de suas diagonais. Consulte a Figura 1.2.

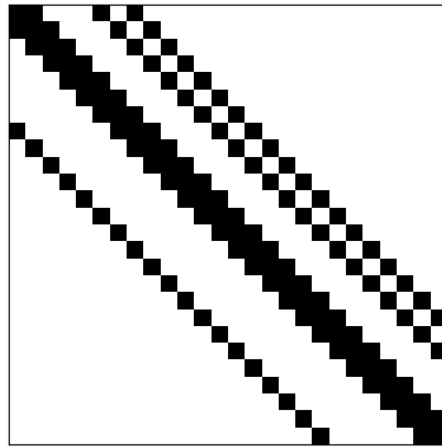


Figura 1.2: Exemplo de uma matriz banda.

Exemplo 1.1.1. Consideremos o seguinte problema de Poisson³

$$-\Delta u = f(x, y), (x, y) \in (0, \pi) \times (0, \pi), \quad (1.25)$$

$$u(0, y) = 0, y \in [0, \pi], \quad (1.26)$$

$$u(\pi, y) = 0, y \in [0, \pi], \quad (1.27)$$

$$u(x, 0) = 0, x \in [0, \pi], \quad (1.28)$$

$$u(x, \pi) = 0, x \in [0, \pi]. \quad (1.29)$$

Para fixarmos as ideias, vamos assumir

$$f(x, y) = \text{sen}(x) \text{sen}(y) \quad (1.30)$$

Vamos empregar o Método de Diferenças Finitas⁴ para computar uma aproximação para a sua solução. Começamos assumindo uma malha uniforme de n^2 nodos

$$x_i = (i - 1)h \quad (1.31)$$

$$y_j = (j - 1)h \quad (1.32)$$

com tamanho de malha $h = \pi/(n - 1)$, $i = 1, 2, \dots, n$ e $j = 1, 2, \dots, n$. Empregando a Fórmula de Diferenças Central⁵ encontramos o seguinte problema discreto associado

$$u_{i,1} = u_{1,j} = 0 \quad (1.33)$$

$$\begin{aligned} & -\frac{1}{h^2}u_{i-1,j} - \frac{1}{h^2}u_{i,j-1} + \frac{4}{h^2}u_{i,j} \\ & -\frac{1}{h^2}u_{i+1,j} - \frac{1}{h^2}u_{i,j+1} = f(x_i, y_j) \end{aligned} \quad (1.34)$$

$$u_{i,n} = u_{n,j} = 0 \quad (1.35)$$

Este é um sistema linear $n^2 \times n^2$. Tomando em conta as condições de contorno, ele pode ser reduzido a um sistema $(n - 2)^2 \times (n - 2)^2$

$$Aw = b \quad (1.36)$$

³Baron Siméon Denis Poisson, 1781 - 1840, matemático, engenheiro e físico francês. Fonte: [Wikipedia](#).

⁴Observamos que $\Delta u = u_{xx} + u_{yy}$.

⁵Consulte mais em [Notas de Aula: Matemática Numérica](#).

usando a enumeração das incógnitas $(i, j) \rightarrow k = i - 1 + (j - 2)(n - 2)$, i.e.

$$u_{i,j} = w_{k=i-1+(j-2)(n-2)}, \quad i, j = 2, \dots, n-2 \quad (1.37)$$

Consulte a Figura 1.3 para uma representação da enumeração em relação a malha.

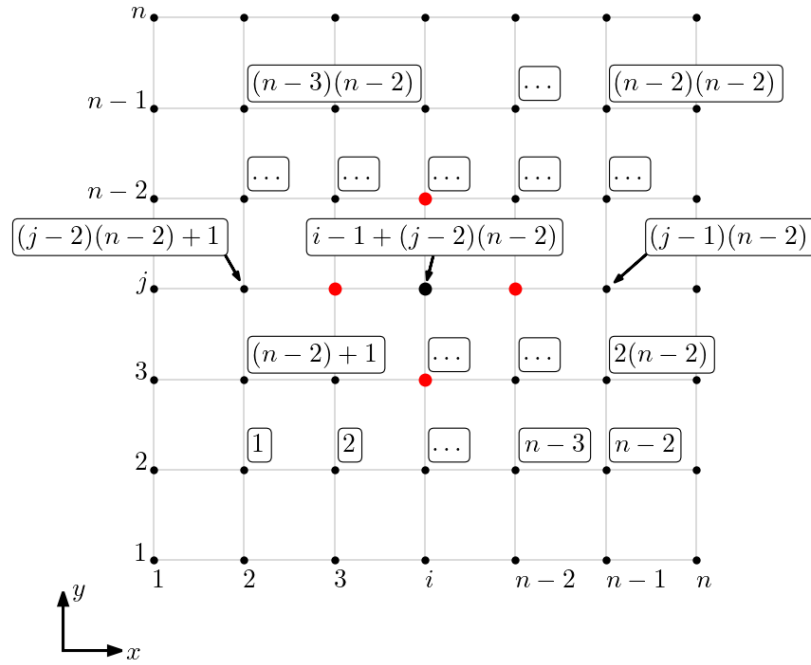


Figura 1.3: Representação da enumeração das incógnitas referente ao problema discutido no Exercício 1.1.1.

Afim de obtermos uma matriz diagonal dominante, vamos ordenar as equações do sistema discreto como segue

- $j = 2, i = 2$:

$$4w_k - w_{k+1} - w_{k+n-2} = h^2 f_{i,j} \quad (1.38)$$

- $j = 2, i = 3, \dots, n-2$:

$$-w_{k-1} + 4w_k - w_{k+1} - w_{k+n-2} = h^2 f_{i,j} \quad (1.39)$$

- $j = 2, i = n - 1$:

$$-w_{k-1} + 4w_k - w_{k+n-2} = h^2 f_{i,j} \quad (1.40)$$

- $j = 3, \dots, n - 2, i = 2$:

$$-w_{k-(n-2)} + 4w_k - w_{k+1} - w_{k+n-2} = h^2 f_{i,j} \quad (1.41)$$

- $j = 3, \dots, n - 2, i = 3, \dots, n - 2$:

$$-w_{k-1} - w_{k-(n-2)} + 4w_k - w_{k+1} - w_{k+n-2} = h^2 f_{i,j} \quad (1.42)$$

- $j = 3, \dots, n - 2, i = n - 1$:

$$-w_{k-1} - w_{k-(n-2)} + 4w_k - w_{k+n-2} = h^2 f_{i,j} \quad (1.43)$$

- $j = n - 1, i = 2$:

$$-w_{k-(n-2)} + 4w_k - w_{k+1} = h^2 f_{i,j} \quad (1.44)$$

- $j = n - 1, i = 3, \dots, n - 2$:

$$-w_{k-1} - w_{k-(n-2)} + 4w_k - w_{k+1} = h^2 f_{i,j} \quad (1.45)$$

- $j = n - 1, i = n - 1$:

$$-w_{k-1} - w_{k-(n-2)} + 4w_k = h^2 f_{i,j} \quad (1.46)$$

Com isso, obtemos uma matriz com 5 bandas, consulte a Figura 1.4.

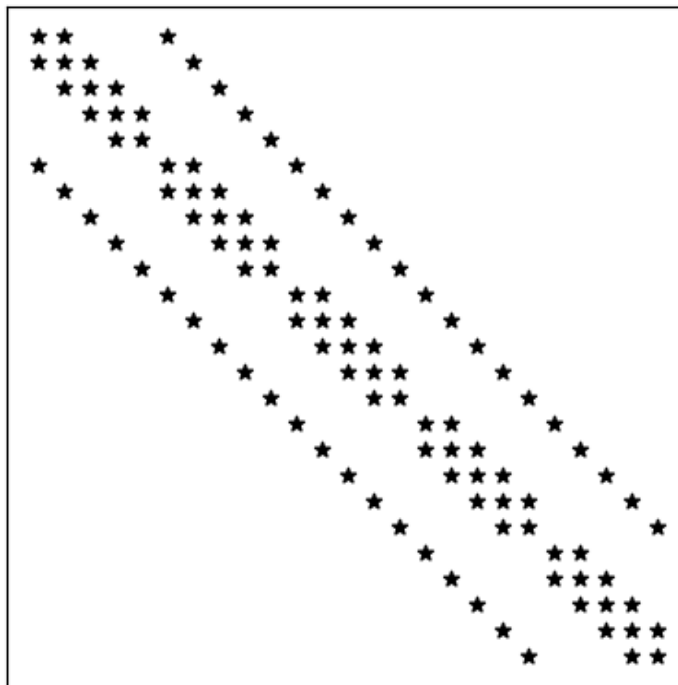


Figura 1.4: Representação da matriz do sistema discreto construído no Exemplo 1.1.1.

Exercício 1.1.3. Consideremos o problema trabalho no Exemplo 1.1.1.

- a) Use a função `scipy.linalg.solve` para computar a solução do problema discreto associado para diferentes tamanhos de malha $h = 10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}$. Compute o erro da solução discreta em relação à solução analítica. Compare as aproximações com a solução analítica

$$u(x, y) = \frac{1}{2} \sin(x) \sin(y). \quad (1.47)$$

- b) Compare a demanda de tempo e memória computacional se a função `scipy.linalg.solve_banded` for empregada na computação da solução discreta.
- c) Baseado no Algoritmo de Thomas, implemente o Método de Eliminação

Gaussiana otimizado para a matriz banda deste problema. Compare com as abordagens dos itens a) e b).

1.1.3 Esquemas de Armazenamento

A ideia é armazenar apenas os elementos não-nulos de uma matriz esparsa, de forma a economizar a demanda de armazenamento computacional. Cuidados devem ser tomados para que a estrutura de armazenamento utilizada seja adequada para a computação das operações matriciais mais comuns.

Formato COO

O formato COO (*COOrdinate format*) é o esquema de armazenamento mais simples. A estrutura de dados consiste em três arranjos: (1) um arranjo contendo as entradas não-nulas da matriz; (2) um arranjo contendo seus índices de linha; (3) um arranjo contendo seus índices de coluna.

Exemplo 1.1.2. Consideremos a seguinte matriz

$$A = \begin{bmatrix} 2. & 0. & 1. & 0. \\ 0. & 3. & 2. & -1. \\ 0 & -1 & -2. & 0. \\ 0 & 0 & 0 & 1. \end{bmatrix} \quad (1.48)$$

O formato COO está disponível na biblioteca SciPy com o método [scipy.sparse.coo_matrix](#)⁶. Neste caso, temos

```
1     import numpy as np
2     import scipy as sp
3     from scipy.sparse import coo_matrix
4
5     data = np.array([2.,1.,3.,2.,-1.,-1.,-2.,1.])
6     row = np.array([0,0,1,1,1,2,2,3])
7     col = np.array([0,2,1,2,3,1,2,3])
8     coo = coo_matrix((data, (row, col)), shape=(4,4))
9     print("coo = \n", coo)
10    print("A = \n", coo.toarray())
```

⁶Versões recentes do SciPy estão migrando para a nomenclatura do NumPy. Consulte mais em [SciPy Sparse](#).

- Vantagens do formato COO são:
 - permite a entrada de dados duplicados (simplicidade);
 - possível conversão rápida entre os formatos CSR e CSC⁷.
- Desvantagens do formato COO são:
 - complexidade em operações aritméticas;
 - complexidade na extração de submatrizes.

Observação 1.1.1. O SciPy conta com vários métodos para o tratamento e operação com matrizes esparsas armazenadas no formato COO. Consulte mais em [scipy.sparse.coo_matrix](https://docs.scipy.org/doc/scipy/reference/sparse.coo_matrix.html).

Formato CSR

O formato *Compressed Sparse Row* (CSR) é uma variação do COO que busca diminuir a alocação de dados repetidos. Assim como o COO, o formato conta com três arranjos d , c , p :

- d é o arranjo contendo os elementos não-nulos da matriz, ordenados por linhas (i.e., da esquerda para direita, de cima para baixo);
- c é o arranjo contendo o índice das colunas das entradas não-nulas da matriz (como no formato COO);
- p é um arranjo cujos elementos são a posição no arranjo c em que cada linha da matriz começa a ser representada. O número de elementos de i -ésima linha da matriz dado por $p_{j+1} - p_j$.

Exemplo 1.1.3. No Exemplo 1.1.2, alocamos a matriz

$$A = \begin{bmatrix} 2. & 0. & 1. & 0. \\ 0. & 3. & 2. & -1. \\ 0 & -1 & -2. & 0. \\ 0 & 0 & 0 & 1. \end{bmatrix} \quad (1.49)$$

no formato COO. Aqui, vamos converter a alocação para o formato CSR e, então, verificar seus atributos.

⁷Estes formatos são mais eficientes para a computação matricial e são apresentados na sequência.


```
1 from scipy.sparse import csr_matrix
2 csr = coo.tocsr()
3 d = csr.data
4 c = csr.indices
5 p = csr.indptr
6 print(d)
7 print(c)
8 print(p)
```

Para fixarmos as ideias, temos

$$d = (2., 1., 3., 2., -1., -1., -2., 1.) \quad (1.50)$$

$$c = (0, 2, 1, 2, 3, 1, 2, 3) \quad (1.51)$$

$$p = (0, 2, 5, 7, 8) \quad (1.52)$$

Assim sendo, o elemento $p[i=2] = 5$ aponta para o $c[k=5] = 1$, o que fornece que $A[i=2, j=1] = d[k] = -1$. Verifique!

- Vantagens do formato CSR:
 - operações aritméticas eficientes;
 - fatiamento por linhas eficiente;
 - multiplicação matriz vetor eficiente.
- Desvantagens do formato CSR:
 - fatiamento por colunas não eficiente;
 - custo elevado de realocamento com alteração da esparsidade da matriz.

Observação 1.1.2. O SciPy conta com vários métodos para o tratamento e operação com matrizes esparsas armazenadas no formato CSR. Consulte mais em [scipy.sparse.csr_matrix](https://docs.scipy.org/doc/scipy/reference/sparse.csr_matrix.html).

Formato CSC

O formato *Compressed Sparse Column* (CSC) é uma variação análoga do CSR, mas para armazenamento por colunas. O formato conta com três arranjos d , l , p :

- d é o arranjo contendo os elementos não-nulos da matriz, ordenados por colunas (i.e., de cima para baixo, da esquerda para direita);
- l é o arranjo contendo o índice das linhas das entradas não-nulas da matriz;
- p é um arranjo cujos elementos são a posição no arranjo l em que cada coluna da matriz começa a ser representada. O número de elementos de j -ésima coluna da matriz dado por $p_{j+1} - p_j$.

Exemplo 1.1.4. No Exemplo 1.1.2, alocamos a matriz

$$A = \begin{bmatrix} 2. & 0. & 1. & 0. \\ 0. & 3. & 2. & -1. \\ 0 & -1 & -2. & 0. \\ 0 & 0 & 0 & 1. \end{bmatrix} \quad (1.53)$$

no formato COO. Aqui, vamos converter a alocação para o formato CSC e, então, verificar seus atributos.

```

1      from scipy.sparse import csc_matrix
2      csc = coo.tocsc()
3      d = csc.data
4      l = csc.indices
5      p = csc.indptr
6      print(d)
7      print(l)
8      print(p)
```

Para fixarmos as ideias, temos

$$d = (2., 3., -1., 1., 2., -2., -1., 1.) \quad (1.54)$$

$$l = (0, 1, 2, 0, 1, 2, 1, 3) \quad (1.55)$$

$$p = (01368) \quad (1.56)$$

Assim sendo, o elemento $p[j=2] = 3$ aponta para o $l[k=3] = 0$, o que informa que $A[i=0, j=2] = d_{\{k\}=1}$. Verifique!

- Vantagens do formato CSC:
 - fatiamento por colunas eficiente;
 - operações aritméticas eficientes;

- multiplicação matriz vetor eficiente⁸.
- Desvantagens do formato CSC:
 - fatiamento por linhas não eficiente;
 - custo elevado de realocamento com alteração da esparsidade da matriz.

Observação 1.1.3. O SciPy conta com vários métodos para o tratamento e operação com matrizes esparsas armazenadas no formato CSC. Consulte mais em [scipy.sparse.csc_matrix](#).

Observação 1.1.4. Além dos formatos COO, CSR e CSC, existem ainda vários outros que podem empregados e que são mais eficientes em determinadas aplicações. Recomendamos a leitura de [3, Seção 3.4] e da documentação do [scipy.sparse](#).

Exercício 1.1.4. Considere o problema de Poisson dado no Exemplo 1.1.1.

- a) Armazene a matriz do problema discreto associado usando o formato COO.
- b) Converta a matriz armazenada para o formato CSR⁹. Então, compute a solução do problema discreto com o método [spsolve](#)¹⁰.
- c) Converta a matriz armazenada para o formato CSC¹¹. Então, compute a solução do problema discreto com o método [spsolve](#).
- d) Compare a eficiência da computação entre os itens b) e c) para tamanhos de malha $h = 10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}$.

⁸CSR é mais eficiente em muitos casos.

⁹Use o método [coo_matrix.tocsr\(\)](#).

¹⁰[scipy.sparse.linalg.spsolve](#) é uma implementação do Método LU otimizado para matrizes esparsas.

¹¹Use o método [coo_matrix.tocsc\(\)](#).

1.2 Métodos Iterativos

1.2.1 GMRES

O GMRES (do inglês, *Generalized Minimal Residual Method*¹²) é um Método de Subespaço de Krylov e é considerado uma das mais eficientes técnicas para a resolução de sistemas lineares gerais e de grande porte (esparsos).

Método de Subespaço de Krylov

A ideia básica é resolver o sistema linear

$$Ax = b \quad (1.57)$$

por um **método de projeção**. Mais especificamente, busca-se uma solução aproximada $x_m \in \mathbb{R}^N$ no subespaço afim $x_0 + \mathcal{K}_m$ de dimensão m , impondo-se a **condição de Petrov¹³-Galerkin¹⁴**

$$b - Ax_m \perp \mathcal{L}_m, \quad (1.58)$$

onde \mathcal{L}_m também é um subespaço de dimensão m . Quando \mathcal{K}_m é um **subespaço de Krylov¹⁵**, i.e.

$$\mathcal{K}_m(A, r_0) = \{r_0, Ar_0, A^2r_0, \dots, A^{m-1}r_0\}, \quad (1.59)$$

temos o Método de Subespaço de Krylov. Aqui, temos o **resíduo**

$$r_0 = b - Ax_0, \quad (1.60)$$

sendo x_0 uma aproximação inicial para a solução do sistema. Notemos que com isso, temos que a aproximação calculada é tal que

$$A^{-1}b \approx x_m = x_0 + q_{m-1}(A)r_0, \quad (1.61)$$

onde q_{m-1} é um dado polinômio de grau $m-1$. No caso particular de $x_0 = 0$, temos

$$A^{-1}b \approx q_{m-1}(A)x_0. \quad (1.62)$$

Diferentes versões deste método são obtidas pelas escolhas do subespaço \mathcal{L}_m e formas de preconditionamento do sistema.

¹²Desenvolvido por Yousef Saad e H. Schultz, 1986. Fonte: [Wikipedia](#).

¹³Georgi Iwanowitsch Petrov, 1912 - 1987, engenheiro soviético. Fonte: [Wikipedia](#).

¹⁴Boris Galerkin, 1871 - 1945, engenheiro e matemático soviético. Fonte: [Wikipédia](#).

¹⁵Alexei Nikolajewitsch Krylov, 1863 - 1945, engenheiro e matemático russo. Fonte: [Wikipédia](#).

GMRES

O GMRES é um Método de Subespaço de Krylov assumindo $\mathcal{L}_m = A\mathcal{K}_m$, com

$$\mathcal{K}_m = \mathcal{K}_m(A, v_1) = \{v_1, Av_1, \dots, A^{m-1}v_1\}, \quad (1.63)$$

onde $v_1 = r_0/\|r_0\|$ é o vetor unitário do resíduo $r_0 = b - Ax_0$ para uma dada aproximação inicial x_0 da solução do sistema $Ax = b$.

Vamos derivar o método observando que qualquer vetor x em $x_0 + \mathcal{K}_m$ pode ser escrito como segue

$$x = x_0 + V_m y \quad (1.64)$$

onde, $V_m = [v_1, \dots, v_m]$ é a matriz $n \times m$ cujas colunas formam uma base ortogonal $\{v_1, \dots, v_m\}$ de \mathcal{K}_m e $y \in \mathbb{R}^m$. Aqui, V_m é computada usando-se o seguinte **Método de Arnoldi¹⁶ - Gram¹⁷ - Schmidt¹⁸ Modificado** [3, Subseção 6.3]:

1. Dado v_1 de norma 1

2. Para $j = 1, \dots, m$:

(a) $w_j := Av_j$

(b) Para $i = 1, \dots, j$:

i. $h_{i,j} := (w_j, v_i)$

ii. $w_j := w_j - h_{i,j}v_i$

(c) $h_{j+1,j} := \|w_j\|$

(d) Se $h_{j+1,j} = 0$, então pare.

(e) $v_{j+1} = w_j/h_{j+1,j}$

Seja, então, $\bar{H}_m = [h_{i,j}]_{i,j=1}^{m+1,m}$ a matriz de Hessenberg¹⁹ cujas entradas não nulas são computadas pelo algoritmo acima (Passos 2(a)i-ii). Pode-se

¹⁶Walter Edwin Arnoldi, 1917 - 1995, engenheiro americano estadunidense. Fonte: [Wikipédia](#).

¹⁷Jørgen Pedersen Gram, 1850 - 1916, matemático dinamarquês. Fonte: [Wikipédia](#).

¹⁸Erhard Schmidt, 1876 - 1959, matemático alemão. Fonte: [Wikipédia](#).

¹⁹Karl Adolf Hessenberg, 1904 - 1959, engenheiro e matemático alemão. Fonte: [Wikipédia](#).

mostrar²⁰ que

$$J(y) = \|b - Ax\| \quad (1.65)$$

$$= \|b - A(x_0 + V_my)\| \quad (1.66)$$

$$= \|\beta e_1 - \tilde{H}_m y\| \quad (1.67)$$

onde, $\beta = \|r_0\|$.

A aproximação GMRES x_m é então computada como

$$x_m = x_0 + V_my_m, \quad (1.68)$$

$$y_m = \min_y \|\beta e_1 - \tilde{H}_m y\| \quad (1.69)$$

Observamos que este último é um pequeno problema de minimização, sendo que requer a solução de um sistema $(m+1) \times m$ de mínimos quadrados, sendo m normalmente pequeno.

Em resumo, a solução GMRES x_m é computada seguindo os seguintes passos:

1. Escolhemos uma aproximação inicial x_0 para a solução de $Ax = b$.
2. Calculamos o resíduo $r_0 = b - Ax_0$.
3. Calculamos o vetor unitário $v_1 = r_0/\|r_0\|$.
4. Usamos o Método de Arnoldi-Gram-Schmidt Modificado para calculamos uma base ortogonal V_m de \mathcal{K}_m e a matriz de Hessenberg \tilde{H}_m associada.
5. Calculamos $y_m = \min_y \|\beta e_1 - \tilde{H}_m y\|$.
6. Calculamos $x_m = x_0 + V_my$.

Observação 1.2.1 (Convergência). Pode-se mostrar que o GMRES converge em ao menos n passos.

Observação 1.2.2 (GMRES com a ortogonalização de Householder). No algoritmo acima, o Método Modificado de Gram-Schmidt é utilizado no processo de Arnoldi. Uma versão numericamente mais eficiente é obtida quando a **Transformação de Householder**²¹ é utilizada. Consulte mais em [3, Subsection 6.5.2].

²⁰Consulte [3, Proposição 6.5].

²¹Alston Scott Householder, 1904 - 1993, matemático americano estadunidense. Fonte: Wikipédia.

Observação 1.2.3 (GMRES com Reinicialização). O *Restarted GMRES* é uma variação do método para sistemas que requerem uma aproximação GMRES x_m com m grande. Nestes casos, o método original pode demandar um custo muito alto de memória computacional. A ideia consiste em assumir m pequeno e, caso não suficiente, recalculer a aproximação GMRES com $x_0 = x_m$. Este algoritmo pode ser descrito como segue.

1. Computamos $r_0 = b - Ax_0$, $\beta = \|r_0\|$ e $v_1 = r_0/\beta$
2. Computamos V_m e \hat{H}_m pelo método de Arnoldi
3. Computamos

$$y_m = \min_y \|\beta e_1 - \hat{H}_m y\| \quad (1.70)$$

$$x_m = x_0 + V_m y_m \quad (1.71)$$

4. Se $\|b - Ax_m\|$ é satisfatória, paramos. Caso contrário, setamos $x_0 := x_m$ e voltamos ao passo 1.

A convergência do *Restarted GMRES* não é garantida para matrizes que não sejam positiva-definidas.

Exercício 1.2.1. Considere o problema discreto do Exercício 1.1.2.

- a) Compute a solução com a implementação *Restarted GMRES*

[`scipy.sparse.linalg.gmres`](#).

- b) Por padrão, o intervalo de iterações entre as inicializações é `restart=20`. Compare o desempenho para diferentes intervalos de reinicialização.
- c) Compare o desempenho entre as abordagens dos itens a) e b) frente a implementação do Método de Eliminação Gaussiana disponível em

[`scipy.sparse.linalg.spsolve`](#).

Exercício 1.2.2. Considere o problema discreto trabalhado no Exemplo 1.1.1.

- a) Compute a solução com a implementação *Restarted GMRES*

[`scipy.sparse.linalg.gmres`](#).

- b) Por padrão, o intervalo de iterações entre as inicializações é `restart=20`. Compare o desempenho para diferentes intervalos de reinicialização.
- c) Compare o desempenho entre as abordagens dos itens a) e b) frente a implementação do Método de Eliminação Gaussiana disponível em [scipy.sparse.linalg.spsolve](https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.linalg.spsolve.html).

Exercício 1.2.3. Faça sua própria implementação do método GMRES. Valide-a e compare-a com as resoluções dos Exercícios 1.2.1 e ??.

1.2.2 Método do Gradiente Conjugado

O Método do Gradiente Conjugado é o uma das mais eficientes técnicas iterativas para a resolução de sistema linear com matriz esparsa, simétrica e definida positiva. Portanto, vamos assumir que o sistema

$$Ax = b \quad (1.72)$$

onde, a A é **simétrica** e **definida positiva**.

O método pode ser derivado a partir do ²² [3, Seção 6.7] ou como uma variação do Método do Gradiente. Este é caminho que será adotado aqui.

Método do Gradiente

A ideia é reformular o sistema $Ax = b$ como um problema de minimização. Vamos começar definindo o funcional

$$J(y) = \frac{1}{2}y^T Ay - y^T b. \quad (1.73)$$

O vetor y que minimiza J é a solução de $Ax = b$. De fato, denotando x a solução de $Ax = b$, temos

$$J(y) = \frac{1}{2}y^T Ay - y^T b + \frac{1}{2}x^T Ax - \frac{1}{2}x^T Ax \quad (1.74)$$

$$= \frac{1}{2}(y - x)^T A(y - x) - \frac{1}{2}x^T Ax \quad (1.75)$$

²²Walter Edwin Arnoldi, 1917 - 1995, engenheiro americano estadunidense. Fonte: [Wikipedia](https://pt.wikipedia.org/wiki/Walter_Edwin_Arnoldi).

O termo é independente de y e, portanto, J é mínimo quando

$$\frac{1}{2}(y-x)^T A(y-x) \quad (1.76)$$

é minimizado. Agora, como A é definida positiva²³, o menor valor deste termo ocorre quando $y-x=0$, i.e. $y=x$.

Observamos, também, que o gradiente de J é

$$\nabla J = Ay - b \quad (1.77)$$

i.e., é o oposto do resíduo $r = b - Ay$. Com isso, temos que $y=x$ é a única escolha tal que $\nabla J = 0$. Ainda, temos que ∇J é o vetor que aponta na direção e sentido de maior crescimento de J . Isso nos motiva a aplicarmos a seguinte iteração²⁴

$$x^{(0)} = \text{aprox. inicial} \quad (1.78)$$

$$x^{(k+1)} = x^{(k)} - \alpha_k \nabla J(x^{(k)}) \quad (1.79)$$

onde, $\alpha_k > 0$ é um escalar que regula o tamanho do passo a cada iteração. Lembrando que $-\nabla J = r$, temos que a iteração é equivalente a

$$x^{(k+1)} = x^{(k)} + \alpha_k r^{(k)}. \quad (1.80)$$

Notamos que $x^{(k+1)}$ é um ponto na reta $\{x^{(k)} + \alpha r^{(k)} : \alpha \in \mathbb{R}\}$ que tem a mesma direção de $\nabla J(x^{(k)})$ e passa pelo ponto $x^{(k)}$. O procedimento de escolher um $\alpha^{(k)}$ entre todos os possíveis, é conhecido como pesquisa linear (em inglês, *line search*).

A cada iteração, queremos escolher α_k de forma que $J(x^{(k+1)}) \leq J(x^{(k)})$. Isso pode ser garantido fazendo a seguinte escolha²⁵

$$J(x^{(k+1)}) = \min_{\alpha \in \mathbb{R}} J(x^{(k)} + \alpha r^{(k)}) \quad (1.81)$$

A fim de resolver este problema de minimização, vamos denotar

$$g(\alpha) = J(x^{(k)} + \alpha r^{(k)}). \quad (1.82)$$

²³ $x^T A x > 0$ para todo $x \neq 0$.

²⁴Iteração do Método da Descida mais Rápida.

²⁵Chamada de pesquisa linear exata. Qualquer outra escolha para α é conhecida como pesquisa linear não exata.

Então, observamos que

$$\begin{aligned} g(\alpha) &= \frac{1}{2} \left(x^{(k)} + \alpha r^{(k)} \right)^T A \left(x^{(k)} + \alpha r^{(k)} \right) - \left(x^{(k)} + \alpha r^{(k)} \right)^T b \\ &= \frac{1}{2} x^{(k)T} A x^{(k)} + \frac{\alpha}{2} x^{(k)T} A r^{(k)} + \frac{\alpha}{2} r^{(k)T} A x^{(k)} \\ &\quad + \frac{\alpha^2}{2} r^{(k)T} A r^{(k)} - x^{(k)T} b - \alpha r^{(k)T} b \end{aligned} \quad (1.83)$$

Agora, usando o fato de A ser simétrica, obtemos

$$g(\alpha) = J \left(x^{(k)} \right) + \alpha r^{(k)T} A x^{(k)} + \frac{\alpha^2}{2} r^{(k)T} A r^{(k)} - \alpha r^{(k)T} b \quad (1.84)$$

$$= J \left(x^{(k)} \right) - \alpha r^{(k)T} r^{(k)} + \frac{\alpha^2}{2} r^{(k)T} A r^{(k)} \quad (1.85)$$

a qual, é uma função quadrática. Seu único mínimo, ocorre quando

$$0 = g'(\alpha) \quad (1.86)$$

$$= -r^{(k)T} r^{(k)} + \alpha r^{(k)T} b. \quad (1.87)$$

Logo, encontramos

$$\alpha = \frac{r^{(k)T} r^{(k)}}{r^{(k)T} A r^{(k)}} \quad (1.88)$$

Com isso, temos a iteração do Método do Gradiente

$$x^{(0)} = \text{aprox. inicial} \quad (1.89)$$

$$x^{(k+1)} = x^{(k)} + \alpha_k r^{(k)}, \quad (1.90)$$

$$\alpha_k = \frac{r^{(k)T} r^{(k)}}{r^{(k)T} A r^{(k)}} \quad (1.91)$$

Observação 1.2.4. Observamos que, a cada iteração, precisamos computar $A r^{(k)}$ (no cálculo de α_k) e $A x^{(k)}$ (no cálculo do resíduo). Essas multiplicações matriz-vetor são os passos computacionais mais custosos do método. Podemos otimizar isso usando o fato de que

$$r^{(k+1)} = r^{(k)} - \alpha A r^{(k)}. \quad (1.92)$$

Em construção ...

Referências Bibliográficas

- [1] J.P. Davis and P. Rabinowitz. *Methods of Numerical Integration*. Academic Press, Inc., San Diego, 2. edition, 1984.
- [2] C.T. Kelley. *Iterative Methods for Optimization*. Society for Industrial and Applied Mathematics, Philadelphia, 1999.
- [3] Y. Saad. *Iterative Methods for Sparse Linear Systems*. Society for Industrial and Applied Mathematics, Philadelphia, 2003.
- [4] J. Stoer. *Approximate Calculation of Multiple Integrals*. Prentice-Hall, Inc., 1971.
- [5] D. Watkins. *Fundamentals of Matrix Computations*. John Wiley, New York, 2002.