

4.6

Linear Algebra in Numpy

线性代数

Numpy 中的重要线性代数计算



我的大脑只是一个接收器。宇宙中有一个核心，我们从中获得知识、力量和灵感。这个核心的秘密我没有深入了解，但我知道它的存在。

My brain is only a receiver, in the Universe there is a core from which we obtain knowledge, strength and inspiration. I have not penetrated into the secrets of this core, but I know that it exists.

—— 尼古拉·特斯拉 (Nikola Tesla) | 发明家、物理学家 | 1856 ~ 1943



NumPy 的 linalg 模块

NumPy 库的 linalg 模块提供了许多用于线性代数计算的函数，包括矩阵分解和向量计算。

以下是一些常见的 linalg 函数：

- ▶ `numpy.linalg.inv()`：计算矩阵的逆。
- ▶ `numpy.linalg.pinv()`：计算矩阵的 Moore-Penrose 伪逆。
- ▶ `numpy.linalg.solve()`：求解线性方程组 $Ax = b$ ，其中 A 是一个矩阵， b 是一个向量。
- ▶ `numpy.linalg.lstsq()`：最小二乘解。

linalg 模块还提供了许多向量计算函数，包括：

- ▶ `numpy.linalg.norm()`：计算向量的范数。
- ▶ `numpy.linalg.dot()`：计算向量的点积。

以下是 linalg 中常用的矩阵分解函数：

- ▶ `numpy.linalg.cholesky()`：计算 Cholesky 分解。
- ▶ `numpy.linalg.eig()`：计算矩阵的特征值和特征向量。
- ▶ `numpy.linalg.svd()`：计算奇异值分解。

这些函数在许多科学计算中都非常有用，例如，在机器学习中，可以使用矩阵分解函数进行降维和特征提取，而向量计算函数则可用于计算距离和相似性度量等。需要注意的是，这些函数都要求输入参数为 NumPy 数组，并返回 NumPy 数组作为输出。



什么是矩阵分解？

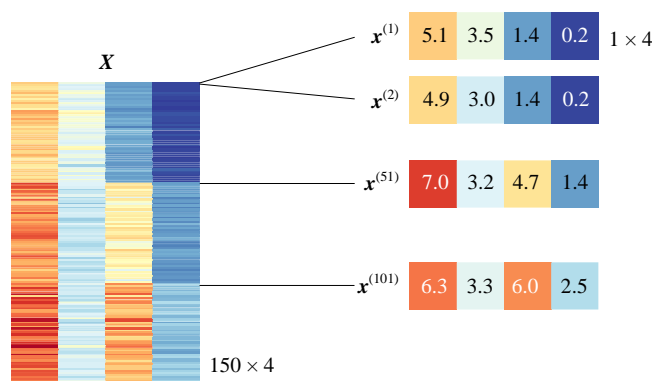
矩阵分解是一种将一个矩阵分解为若干个矩阵的乘积的数学技术。这种分解可以帮助我们更好地理解和处理矩阵数据。常见的矩阵分解包括 Cholesky 分解、特征值分解 (EVD)、奇异值分解 (SVD) 等等。矩阵分解在很多领域都有广泛的应用，比如在机器学习、数据分析、信号处理、图像处理等方面。



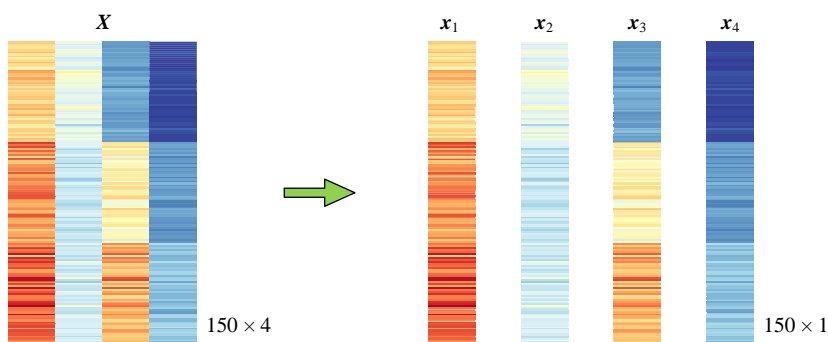
本节配套的 Jupyter Notebook 文件是 BK_2_Topic_4.06_1.ipynb。

拆解矩阵

本书前文提到鸢尾花数据矩阵 X 的形状为 150×4 。也就是说，如图 1 热图所示， X 可以看成是由 150 个行向量上下堆叠而成。每个行向量的形状为 1×4 。图 1 特别展示了 $\mathbf{x}^{(1)}$ (X 第 1 行，数组第 0 行)、 $\mathbf{x}^{(2)}$ (X 第 2 行，数组第 1 行)、 $\mathbf{x}^{(51)}$ (X 第 51 行，数组第 50 行)、 $\mathbf{x}^{(101)}$ (X 第 101 行，数组第 100 行)。

图 1. X 可以看做由一组行向量构成

此外，如图 1 热图所示， X 可以看成是由 4 个列向量左右排列而成，即 $X = [x_1, x_2, x_3, x_4]$ 。每个列向量的形状为 150×1 。请大家回忆，我们如何设定索引获得 NumPy 数组的行向量、列向量。

图 2. X 可以看做由一组列向量构成

几何角度看向量

在二维空间中，一个向量 a 可以表示为一个有序的数对 (a_1, a_2) 、 $[a_1, a_2]$ 、 $[a_1, a_2]^T$ 。向量也可以用一个有向线段来表示，线段的起点为原点 $(0, 0)$ ，终点为 (a_1, a_2) 。其中， a_1 表示向量在水平方向上的投影； a_2 表示向量纵轴方向上的投影。

用勾股定理，我们可求得图 3 中向量 a 的长度，即向量的模，为 $\|a\| = \sqrt{a_1^2 + a_2^2}$ 。在 NumPy 中计算向量模的函数为 `numpy.linalg.norm()`。

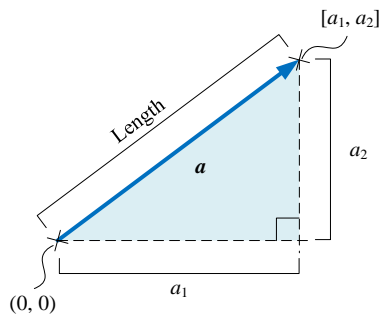


图 3. 向量起点、终点、大小和方向

在 BK_2_Topic_4.06_1.ipynb 中，我们还计算了 $\mathbf{x}^{(1)}$ 、 $\mathbf{x}^{(2)}$ 、 $\mathbf{x}^{(51)}$ 、 $\mathbf{x}^{(101)}$ 这四个向量的单位向量。单位向量是长度为 1 的向量，可以用来表示某个向量方向。比如，的单位向量就是 $\mathbf{x}^{(1)}$ 除以自己的模 $\|\mathbf{x}^{(1)}\|$ ，即 $\mathbf{x}^{(1)} / \|\mathbf{x}^{(1)}\|$ 。



什么是向量的模？

向量的模 (也称为向量的长度) 是指一个向量从原点到其终点的距离，它是一个标量，表示向量的大小。向量的模通常用两个竖线 $\|\mathbf{a}\|$ 来表示，其中 \mathbf{a} 表示向量。对于 n 维向量 $\mathbf{a} = [a_1, a_2, \dots, a_n]$ ，它的模定义为 $\|\mathbf{a}\| = \sqrt{a_1^2 + a_2^2 + \dots + a_n^2}$ 。 $\|\mathbf{a}\|$ 就是向量各个分量的平方和的平方根。这个公式可以用勾股定理推导得出，因为一个向量的模就是从原点到它的终点的距离，而这个距离可以用勾股定理计算。比如，2 维向量 $\mathbf{a} = [3, 4]$ 的模 (长度) 为 $\|\mathbf{a}\| = \sqrt{3^2 + 4^2} = 5$ 。

向量内积

本书前文在讲 for 循环时介绍过 **向量内积** (inner product)，又叫 **标量积** (scalar product)、**点积** (dot product)。给定两个 n 维向量 $\mathbf{a} = [a_1, a_2, \dots, a_n]$ 和 $\mathbf{b} = [b_1, b_2, \dots, b_n]$ ，它们的内积定义为 $\mathbf{a} \cdot \mathbf{b} = a_1b_1 + a_2b_2 + \dots + a_nb_n$ 。内积结果 $\mathbf{a} \cdot \mathbf{b}$ 显然为标量。

如图 4 所示，我们分别计算向量内积 $\mathbf{x}^{(1)} \cdot \mathbf{x}^{(2)}$ 、 $\mathbf{x}^{(1)} \cdot \mathbf{x}^{(51)}$ 、 $\mathbf{x}^{(1)} \cdot \mathbf{x}^{(101)}$ 。建议大家在 JupyterLab 中用手输入算式计算图中三个向量内积。

再次强调，向量内积的运算前提是两个向量维数相同，结果为标量。NumPy 中计算向量内积的函数为 `numpy.dot()`。

$$\begin{array}{ccc}
 \mathbf{x}^{(1)} & \cdot & \mathbf{x}^{(2)} \\
 \begin{bmatrix} 5.1 & 3.5 & 1.4 & 0.2 \end{bmatrix} & \cdot & \begin{bmatrix} 4.9 & 3.0 & 1.4 & 0.2 \end{bmatrix} & = & \mathbf{x}^{(1)} \cdot \mathbf{x}^{(2)} \\
 & & & & 37.49 \\
 \\
 \mathbf{x}^{(1)} & \cdot & \mathbf{x}^{(51)} \\
 \begin{bmatrix} 5.1 & 3.5 & 1.4 & 0.2 \end{bmatrix} & \cdot & \begin{bmatrix} 7.0 & 3.2 & 4.7 & 1.4 \end{bmatrix} & = & \mathbf{x}^{(1)} \cdot \mathbf{x}^{(51)} \\
 & & & & 53.76 \\
 \\
 \mathbf{x}^{(1)} & \cdot & \mathbf{x}^{(101)} \\
 \begin{bmatrix} 5.1 & 3.5 & 1.4 & 0.2 \end{bmatrix} & \cdot & \begin{bmatrix} 6.3 & 3.3 & 6.0 & 2.5 \end{bmatrix} & = & \mathbf{x}^{(1)} \cdot \mathbf{x}^{(101)} \\
 & & & & 52.58
 \end{array}$$

图 4. 向量内积

向量夹角

在 BK_2_Topic_4.06_1.ipynb 中，我们计算得到 $\mathbf{x}^{(1)}$ 、 $\mathbf{x}^{(2)}$ 的夹角约为 3° ， $\mathbf{x}^{(1)}$ 、 $\mathbf{x}^{(51)}$ 的夹角约为 22° ， $\mathbf{x}^{(1)}$ 、 $\mathbf{x}^{(101)}$ 的夹角约为 31° 。

这显然不是巧合， $\mathbf{x}^{(1)}$ 、 $\mathbf{x}^{(2)}$ 分别代表两朵鸢尾花，它们同属 Setosa，因此最为相似。而 $\mathbf{x}^{(51)}$ 属于 Versicolour， $\mathbf{x}^{(101)}$ 属于 Virginica。这就是向量夹角在机器学习中的一个应用举例。



什么是向量夹角？

向量夹角是指两个向量之间的夹角，它是一个标量，通常用弧度或角度来表示。向量夹角的计算是通过向量内积和向量模的关系得出的。对于两个非零向量 \mathbf{a} 和 \mathbf{b} ，它们的夹角 θ 定义为 $\cos(\theta) = (\mathbf{a} \cdot \mathbf{b}) / (\|\mathbf{a}\| \|\mathbf{b}\|)$ 。其中 $\mathbf{a} \cdot \mathbf{b}$ 是向量 \mathbf{a} 和 \mathbf{b} 的内积， $\|\mathbf{a}\|$ 和 $\|\mathbf{b}\|$ 分别是向量 \mathbf{a} 和 \mathbf{b} 的模。注意，这个公式只适用于非零向量，因为对于零向量，它没有方向，因此无法定义夹角。此外， $\cos(\theta)$ 可以看成是 \mathbf{a} 和 \mathbf{b} 的单位向量的向量内积，即 $\cos(\theta) = (\mathbf{a}/\|\mathbf{a}\|) \cdot (\mathbf{b}/\|\mathbf{b}\|)$ 。

通过向量夹角的计算，我们可以判断两个向量之间的相对方向。如果两个向量的夹角为零度，表示它们的方向相同；如果夹角为 90 度，表示它们互相垂直；如果夹角为 180 度，表示它们的方向相反。在机器学习中，可以通过计算向量夹角来度量两个样本之间的相似性。

矩阵乘法

本书前文介绍过矩阵乘法，假设 \mathbf{A} 是一个 $m \times n$ 的矩阵， \mathbf{B} 是一个 $n \times p$ 的矩阵，则它们的乘积 $\mathbf{C} = \mathbf{AB}$ 是一个 $m \times p$ 的矩阵，相当于“消去” n 。在 NumPy 中矩阵乘法的运算符为 $@$ 。

在本节配套的 Jupyter Notebook 文件中大家可以看到两个有趣的矩阵乘法。

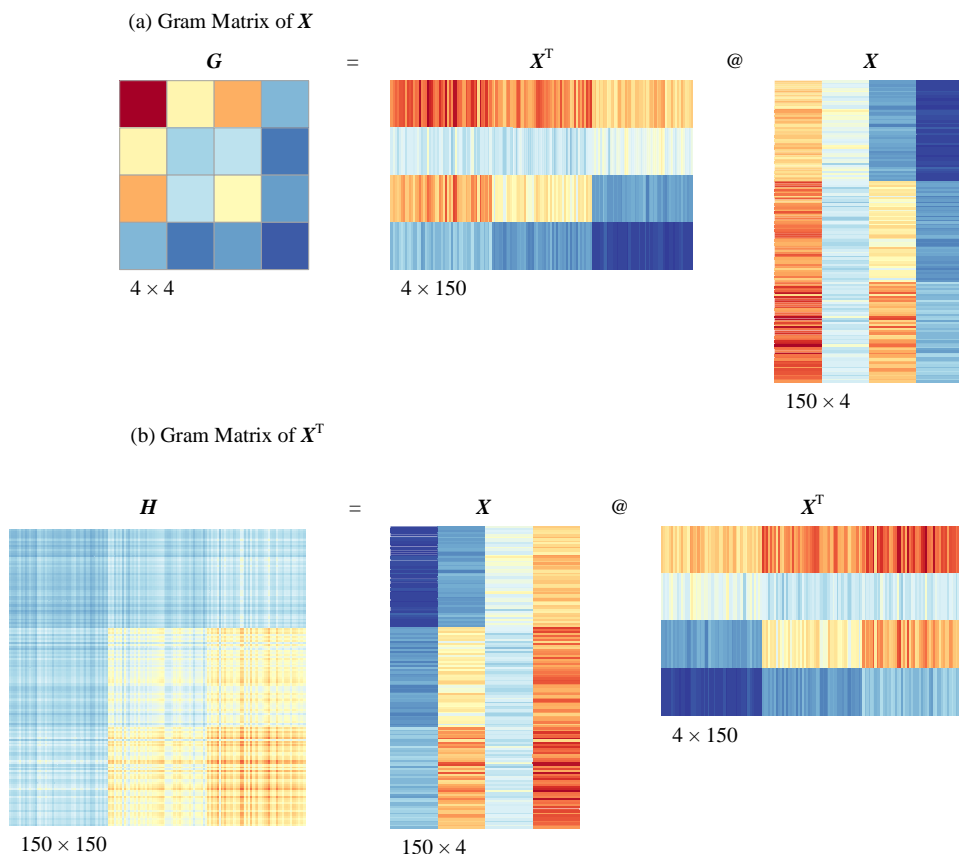


图 5. 两个格拉姆矩阵

如图 5 (a) 所示, 鸢尾花数据矩阵的转置 X^T 乘 X 得到 G 。 X^T 的形状为 4×150 , X 的形状为 150×4 。 $G = X^T X$ 的结果形状为 4×4 。 G 有自己的名字, 叫 X 的格拉姆矩阵 (Gram matrix)。图 5 (b) 所示的 $H = XX^T$ 的结果形状为 150×150 。 H 相当是 X^T 的格拉姆矩阵。



什么是格拉姆矩阵?

格拉姆矩阵 (Gram matrix) 是一个重要的矩阵, 它由向量集合的内积组成。给定一个向量集合 $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, 则其对应的格拉姆矩阵 G 定义为 $G = [g_{ij}]$, 其中 $g_{ij} = \mathbf{x}_i \cdot \mathbf{x}_j$, 表示第 i 个向量和第 j 个向量的内积。格拉姆矩阵是对称矩阵。

格拉姆矩阵在许多应用中都有广泛的应用, 例如在机器学习中的支持向量机 (Support Vector Machine, SVM) 算法和核方法 (kernel method) 中, 格拉姆矩阵可以用来计算向量之间的相似度和距离, 从而实现非线性分类和回归。此外, 格拉姆矩阵也可以用于矩阵分解、图像处理、信号处理等领域。

格拉姆矩阵有很多有趣的性质, 《矩阵力量》一册将详细介绍。这里大家仅仅需要知道格拉姆矩阵为对称矩阵。 G 的主对角线上元素是 $\mathbf{x}_i^T \mathbf{x}_i$, 即 $\mathbf{x}_i \cdot \mathbf{x}_i$ 。如图 6 上图所示, G 的主对角线第一元素 $g_{1,1} = \mathbf{x}_1^T \mathbf{x}_1 = \mathbf{x}_1 \cdot \mathbf{x}_1$ 。如图 6 下图所示, G 的主对角线第二元素 $g_{2,2} = \mathbf{x}_2^T \mathbf{x}_2 = \mathbf{x}_2 \cdot \mathbf{x}_2$ 。请大家自行计算 G 的主对角线剩余两个元素。

如图 7 所示, 显然 $g_{2,1} = g_{1,2}$ 。也就是说, $\mathbf{x}_2^T \mathbf{x}_1 = \mathbf{x}_1^T \mathbf{x}_2 = \mathbf{x}_2 \cdot \mathbf{x}_1 = \mathbf{x}_1 \cdot \mathbf{x}_2$ 。

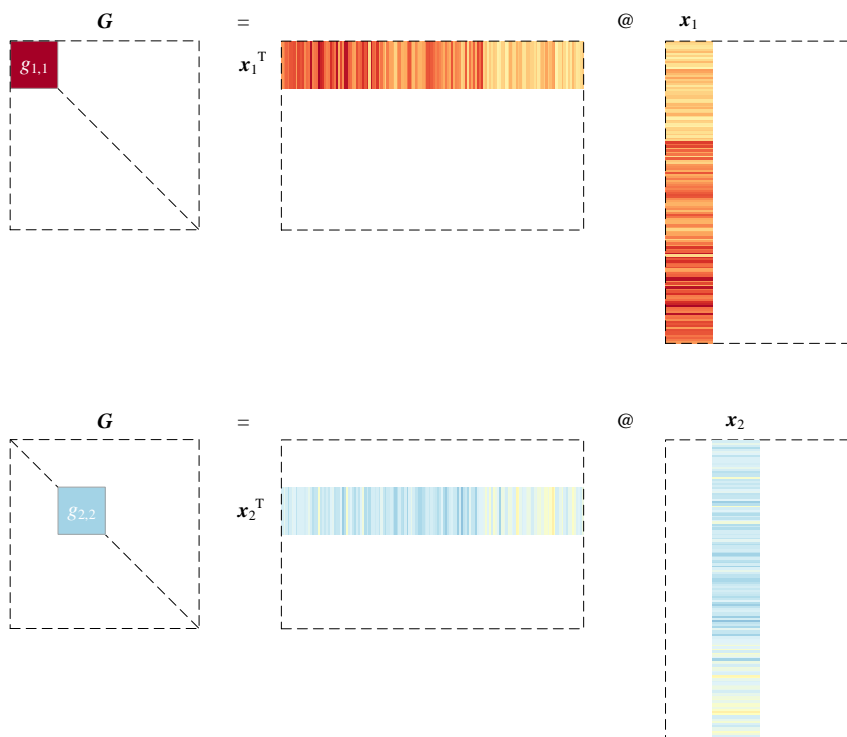


图 6. 格拉姆矩阵 G 主对角线元素

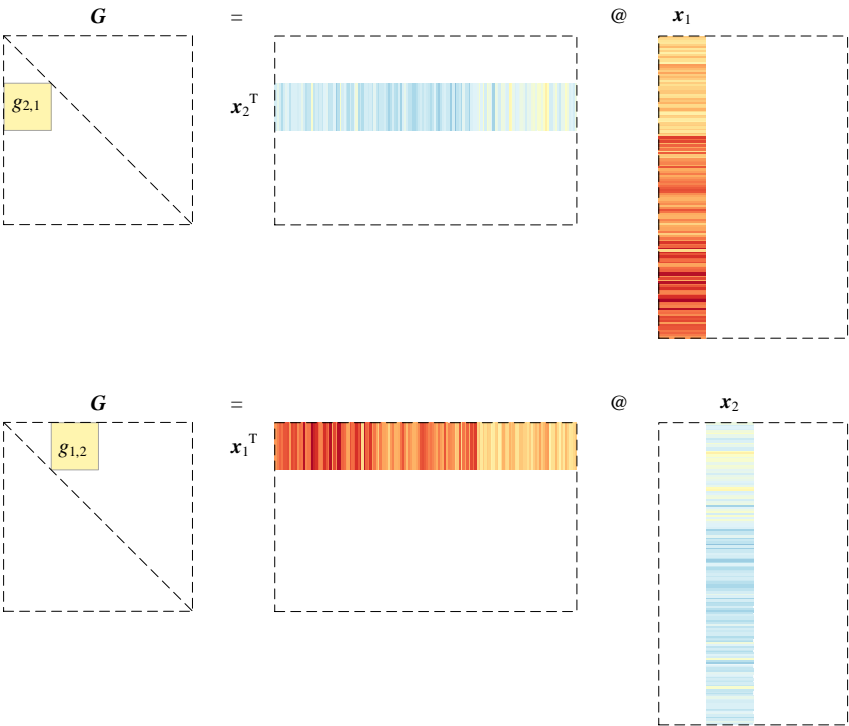


图 7. G 为对称矩阵

矩阵的逆

矩阵的逆可以被看作是一种倒数的概念。并不是所有格拉姆矩阵，恰好前文的格拉姆矩阵 G 存在逆，记做 G^{-1} 。如图 8 所示， G 乘 G^{-1} 结果为单位阵 I 。不难看出， G^{-1} 也是个对称矩阵。

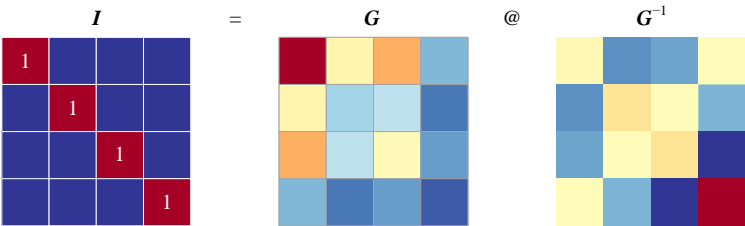


图 8. 格拉姆矩阵 G 的逆



什么是矩阵的逆？

矩阵的逆是一个重要的概念，它是指对于一个可逆的 (即非奇异的) $n \times n$ 矩阵 A ，存在一个 $n \times n$ 矩阵 B ，使得 $AB = BA = I$ ，其中 I 是单位矩阵。 B 被称为 A 的逆矩阵，通常用 A^{-1} 表示。矩阵的逆可以被看作是一种倒数的概念，它可以使我们在矩阵运算中除以矩阵，从而解决线性方程组和其他问题。如果我们需要求解一个线性方程组 $Ax = b$ ，其中 A 是一个可逆矩阵，那么可以使用矩阵的逆来计算 $x = A^{-1}b$ ，从而得到方程的解。需要注意的是，并非所有矩阵都有逆矩阵，只有可逆矩阵才有逆矩阵。对于一个不可逆矩阵，它可能是奇异的 (即行列式为 0)，也可能是非方阵。在实际应用中，矩阵的逆通常通过 LU 分解、QR 分解、Cholesky 分解等方法来计算，而不是直接求解逆矩阵。

Cholesky 分解

所幸前文的格拉姆矩阵 G 也是个正定矩阵 (positive definite matrix)，我们可以对它进行 Cholesky 分解。如图 9 所示， L 是个下三角矩阵，它的转置 L^T 为上三角矩阵。 L 和 L^T 的乘积也相当于“平方”。NumPy 中完成 Cholesky 分解的函数为 `numpy.linalg.cholesky()`。

《矩阵力量》第 12 章专门讲解 Cholesky 分解，这本书第 21 章将介绍正定性。

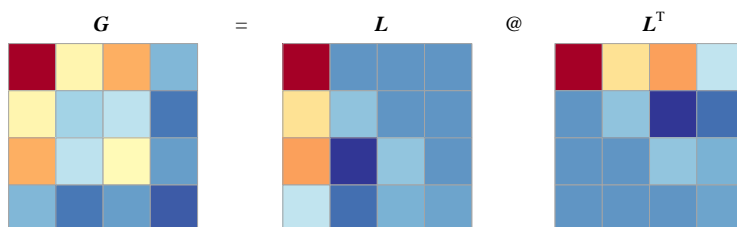


图 9. 对格拉姆矩阵 G 进行 Cholesky 分解



什么是 Cholesky 分解？

Cholesky 分解是一种将对称正定矩阵分解为下三角矩阵和其转置矩阵乘积的数学技术。给定一个对称正定矩阵 A ，Cholesky 分解可以将其表示为 $A = LL^T$ ，其中 L 是下三角矩阵， L^T 是其转置矩阵。Cholesky 分解是一种高效的矩阵分解方法，它可以在数值计算中减少误差，同时可以加速线性方程组的求解，特别是对于大型的稠密矩阵。因此，Cholesky 分解在很多领域都有广泛的应用，例如统计学、金融学、物理学、工程学等。Cholesky 分解也是一些高级技术的基础，例如蒙特卡洛模拟、Kalman 滤波等等。



什么是正定矩阵？

Cholesky 分解是一种将对称正定矩阵分解为下三角矩阵和其转置矩阵乘积的数学技术。给定一个对称正定矩阵 A ，Cholesky 分解可以将其表示为 $A = LL^T$ ，其中 L 是下三角矩阵， L^T 是其转置矩阵。Cholesky 分解是一种高效的矩阵分解方法，它可以在数值计算中减少误差，同时可以加速线性方程组的求解，特别是对于大型的稠密矩阵。因此，Cholesky 分解在很多领域都有广泛的应用，例如统计学、金融学、物理学、工程学等。Cholesky 分解也是一些高级技术的基础，例如蒙特卡洛模拟、Kalman 滤波等等。

特征值分解 EVD

图 10 所示为对格拉姆矩阵 G 的特征值分解。 V 的每一列对应特征向量， A 的主对角线元素为特征值。

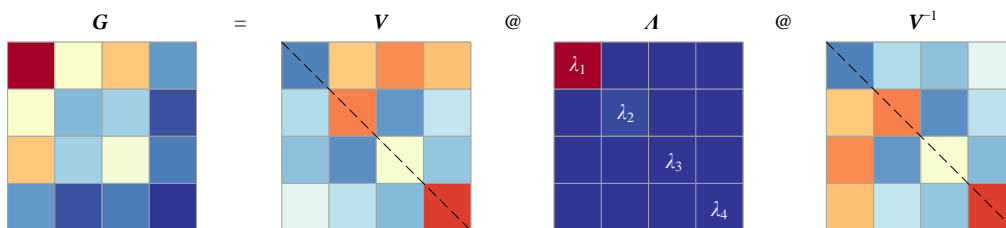


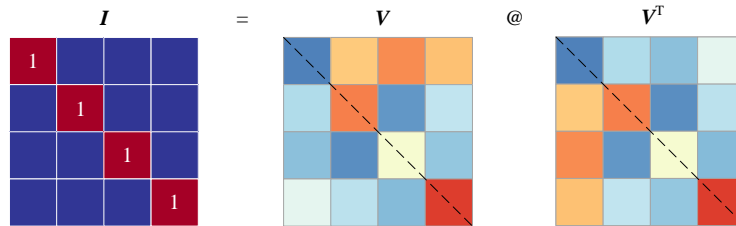
图 10. 对格拉姆矩阵 G 进行 EVD 分解



什么是特征值分解？

特征值分解 (Eigenvalue Decomposition, EVD) 是一种将一个方阵分解为一组特征向量和特征值的数学技术。对于一个 $n \times n$ 的矩阵 A ，如果存在非零向量 v 和常数 λ ，使得 $Av = \lambda v$ ，那么 v 就是矩阵 A 的特征向量， λ 就是对应的特征值。将所有特征向量排列成一个矩阵 V ，将所有特征值排列成一个对角方阵 Λ ，那么矩阵 A 就可以表示为 $A = V\Lambda V^{-1}$ 。特征值分解可以帮助我们理解矩阵的性质和结构，以及实现很多数学算法。它在很多领域都有广泛的应用，比如图像处理、机器学习、信号处理、量子力学等。特征值分解也是一些高级技术的基础，例如奇异值分解、QR 分解、LU 分解等。

仔细观察，大家可以已经发现图 10 中 V 和 V^{-1} 关于主对角线对称，即 $V^T = V^{-1}$ 。这并不是巧合，原因是格拉姆矩阵 G 为对称矩阵。而对称矩阵的特征值分解又叫谱分解 (spectral decomposition)。也就是说， G 的谱分解可以写成 $G = V\Lambda V^T$ 。

图 11. 谱分解中 V 的特点

《矩阵力量》第 13、14 章专门讲解特征值分解。

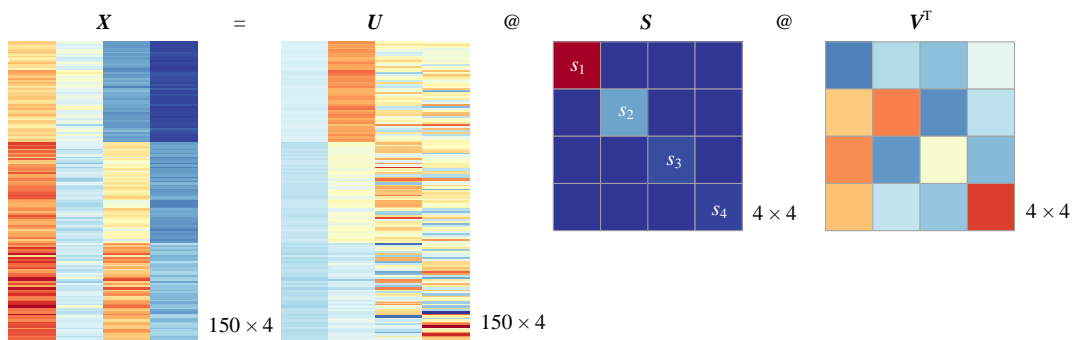


什么是谱分解?

谱分解 (Spectral Decomposition) 是将对称矩阵分解为一组特征向量和特征值的数学技术，即对称矩阵的特征值分解。对于一个对称矩阵 A ，谱分解可以将其分解为 $A = Q\Lambda Q^T$ ，其中 Q 是由矩阵 A 的特征向量组成的正交矩阵， Λ 是由矩阵 A 的特征值组成的对角矩阵。谱分解在很多领域都有广泛的应用，例如图像处理、信号处理、量子力学等。谱分解可以帮助我们理解对称矩阵的性质和结构，从而帮助我们分析和处理各种问题。谱分解也是很多高级技术的基础，例如奇异值分解、主成分分析、矩阵函数等。

奇异值分解 SVD

奇异值分解可谓“最重要的矩阵分解，没有之一”。图 12 所示为对鸢尾花数据矩阵 X 的奇异值分解。图 12 中 S 的主对角线上的元素叫奇异值。大家会在 BK_2_Topic_4.06_1.ipynb 看到，图 10 中特征值开方的结果就是图 12 中的奇异值，这当然不是巧合！

图 12. 对 X 进行 SVD 分解

什么是奇异值分解?

奇异值分解 (Singular Value Decomposition, SVD) 是一种将一个矩阵分解为三个矩阵乘积的数学技术。给定一个矩阵 A ，它可以表示为 $A = U\Sigma V^T$ ，其中 U 和 V 是正交矩阵， Σ 是对角矩阵，对角线上的元素称为奇异值。SVD 可以将一个矩阵的信

息分解为不同奇异值所对应的向量空间，并按照奇异值大小的顺序进行排序，使得我们可以仅使用前面的奇异值和相应的向量空间来近似地表示原始矩阵。这种分解在降维、压缩、数据分析和模型简化等领域中有着广泛的应用，例如推荐系统、图像压缩、语音识别等。

如图 13 所示， U 的转置 U^T 和自己乘积为单位阵。如图 14 所示， V 和自身转置 V^T 乘积为单位阵。大家是否已经发现图 11 和图 14 竟然相同，这当然也不是巧合。

实际上，图 12 是四种奇异值分解中的一种。《矩阵力量》第 15、16 章专门讲解奇异值分解，并揭开各种“巧合”背后的数学原理。

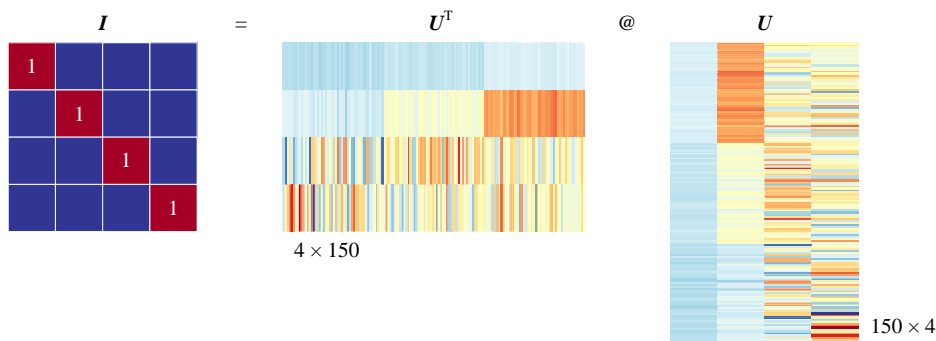


图 13. U 的特点

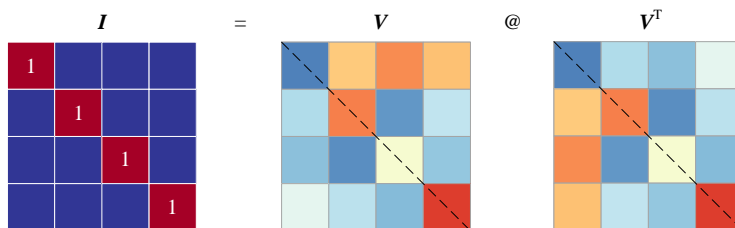


图 14. V 的特点



请大家完成下面 5 道题目。

Q1. 本节配套笔记计算了鸢尾花数据矩阵 X 的若干行向量的模、单位向量、夹角，请大家计算 X 的 4 个列向量的模、单位向量、两两列向量内积、两两夹角。并说明两两列向量内积和图 5 (a) 中格拉姆矩阵的关系。

Q2. 请大家用热图可视化图 5 (a) 中的 G 的第 2 行第 3 列元素如何计算得到。

Q4. 请对图 5 (b) 中的格拉姆矩阵进行 Cholesky 分解，并搜索报错的原因。

Q4. 请对图 5 (b) 中的格拉姆矩阵进行特征值分解，并比较其特征值和图 10 中特征值关系。

Q5. 请对 X^T 进行奇异值分解，比较和图 12 中 SVD 分解的关系。

* 本节不提供答案。