

5

Dive into Matrix Multiplication

矩阵乘法

代数、几何、统计、数据交融的盛宴



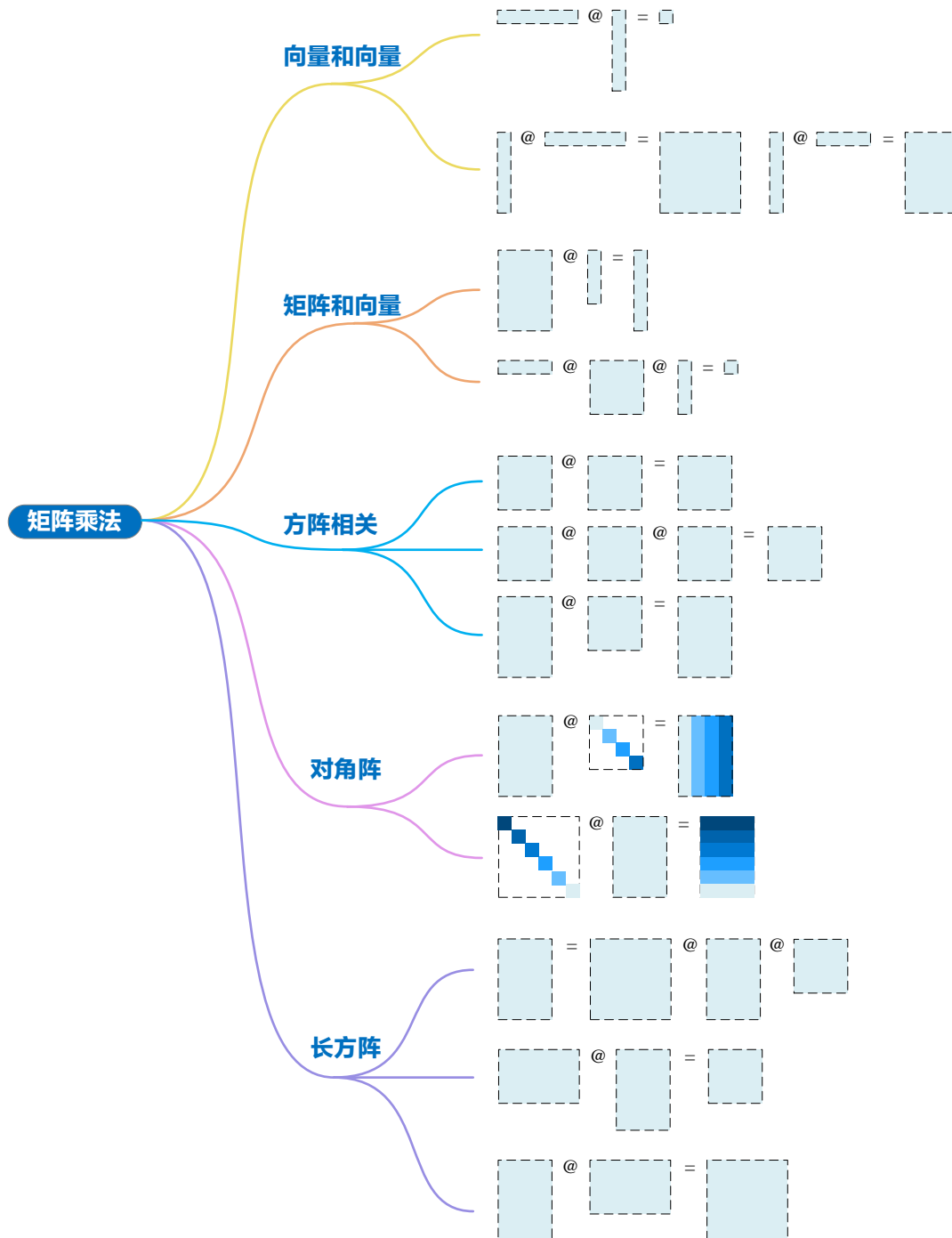
只要持续进步，千万别泼冷水，哪怕蜗行牛步。

Never discourage anyone who continually makes progress, no matter how slow.

—— 柏拉图 (Plato) | 古希腊哲学家 | 424/423 ~ 348/347 BC



- numpy.array() 构造多维矩阵/数组
- numpy.einsum() 爱因斯坦求和约定
- numpy.linalg.inv() 求矩阵逆
- numpy.matrix() 构造二维矩阵
- numpy.multiply() 矩阵逐项积
- numpy.random.random_integers() 生成随机整数
- seaborn.heatmap() 绘制数据热图



本 PDF 文件为作者草稿，发布目的为方便读者在移动终端学习，终稿内容以清华大学出版社纸质出版物为准。

版权归清华大学出版社所有，请勿商用，引用请注明出处。

代码及 PDF 文件下载：<https://github.com/Visualize-ML>

本书配套微课视频均发布在 B 站——生姜 DrGinger：<https://space.bilibili.com/513194466>

欢迎大家批评指教，本书专属邮箱：jiang.visualize.ml@gmail.com

5.1 矩阵乘法：形态丰富多样

矩阵乘法是线性映射的灵魂。因此，矩阵乘法是矩阵运算中最重要的规则，没有之一！

矩阵乘法的规则本身并不难理解；但是，横在我们面前最大的困难是——矩阵乘法的灵活性。这种灵活性主要体现在矩阵乘法不同视角、矩阵乘法形态的多样性这两方面。

本书前文和大家讨论了矩阵乘法的两个视角，本书后续还将在分块矩阵中继续探讨矩阵乘法更多视角。而本章将介绍常见矩阵乘法形态。

▲ 注意，学习本章时，请大家多从代数、几何、数据、统计几个角度理解不同矩阵乘法形态，特别是几何和数据这两个角度。

本章的作用就是鸟瞰全景，让大家开开眼界，不需要大家关注运算细节。如果你之前曾经系统学过线性代数，这一章会让你有寻他千百度、蓦然回首的感觉！作者在学习线性代数的时候，就特别希望能找到一本书能够把常见的矩阵乘法形态和应用场景都娓娓道来。

如果你刚刚接触线性代数相关内容，不要被本章大量术语吓到，大家现在不需要记住它们。本章可以视作全书重要知识点的总结。希望大家在本书不同学习阶段时，能够不断回头翻阅本章，让自己对矩阵乘法的认识一步步加深。

下面，我们就开始“鸟瞰”各种形态的矩阵乘法。

5.2 向量和向量

给定两个等行数列向量 \mathbf{x} 和 \mathbf{y} ：

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad (1)$$

\mathbf{x} 和 \mathbf{y} 向量内积可以写成 \mathbf{x} 转置乘 \mathbf{y} ，或者 \mathbf{y} 转置乘 \mathbf{x} ：

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{y} \cdot \mathbf{x} = \langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = (\mathbf{x}^T \mathbf{y})^T = \mathbf{y}^T \mathbf{x} = x_1 y_1 + x_2 y_2 + \cdots + x_n y_n = \sum_{i=1}^n x_i y_i \quad (2)$$

(2) 告诉我们， $\mathbf{x}^T \mathbf{y}$ 和 $\mathbf{y}^T \mathbf{x}$ 相当于向量元素分别相乘，再求和，结果为标量。这和向量内积的运算结果完全一致，因此我们常用矩阵乘法替代向量内积运算。

观察图 1， $\mathbf{x}^T \mathbf{y}$ 和 $\mathbf{y}^T \mathbf{x}$ 结果均为标量，相当于 1×1 矩阵；这就是为什么 $\mathbf{x}^T \mathbf{y} = (\mathbf{x}^T \mathbf{y})^T$ 。

如果 \mathbf{x} 和 \mathbf{y} **正交** (orthogonal)，则两者向量内积为 0：

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{y} \cdot \mathbf{x} = \langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = (\mathbf{x}^T \mathbf{y})^T = \mathbf{y}^T \mathbf{x} = 0 \quad (3)$$

➔ 正交相当于“垂直”的推广。本书中出现“正交”最多的场合就是“正交投影 (orthogonal projection)”。本书第 9、10 两章专门讲解“正交投影”。

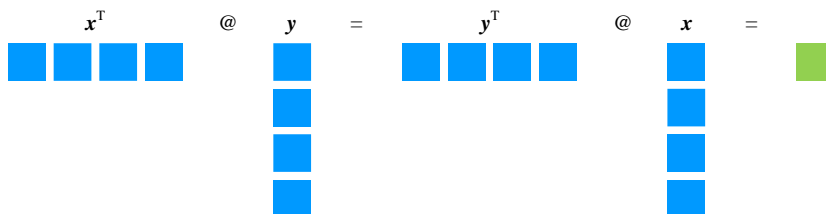


图 1. 标量积

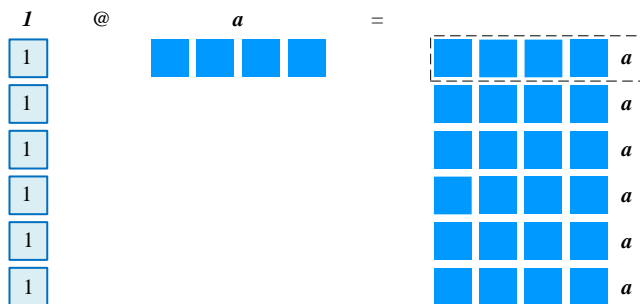
全 1 列向量

全 1 列向量 $\mathbf{1}$ 是非常神奇的存在，多元统计离不开全 1 列向量。下面举几个例子。

如图 2 所示，全 1 列向量 $\mathbf{1}$ 乘行向量 \mathbf{a} ，相当于对行向量 \mathbf{a} 进行复制、向下叠放。 $\mathbf{1} @ \mathbf{a}$ 结果如下：

$$\mathbf{1} @ \mathbf{a} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1} @ \mathbf{a}_{1 \times m} = \begin{bmatrix} \mathbf{a} \\ \mathbf{a} \\ \vdots \\ \mathbf{a} \end{bmatrix}_{n \times m} \quad (4)$$

上式结果为矩阵。复制的份数取决于全 1 列向量 $\mathbf{1}$ 的元素个数。再次强调，上式中 \mathbf{a} 为行向量。

图 2. 复制行向量 \mathbf{a}

类似地，如图 3 所示，列向量 \mathbf{b} 乘全 1 列向量 $\mathbf{1}$ 转置，相当于对列向量 \mathbf{b} 复制、左右排列：

$$\mathbf{b} @ \mathbf{I}^T = \mathbf{b} @ \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}^T = [\mathbf{b} \quad \mathbf{b} \quad \cdots \quad \mathbf{b}] \quad (5)$$

上式结果为矩阵。

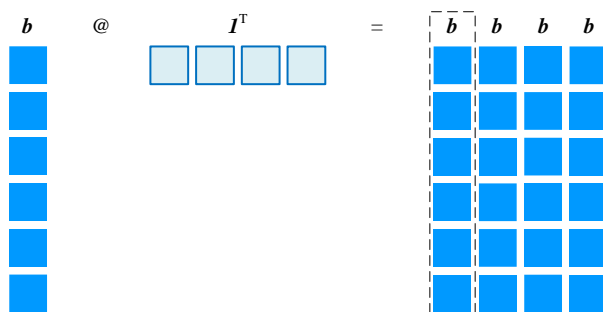


图 3. 复制列向量 \mathbf{b}

统计视角

下式为利用 \mathbf{I} 对列向量 \mathbf{x} 元素求和：

$$\mathbf{I} \cdot \mathbf{x} = \mathbf{I}^T \mathbf{x} = \mathbf{x}^T \mathbf{I} = x_1 + x_2 + \cdots x_n = \sum_{i=1}^n x_i \quad (6)$$

上式结果为标量。

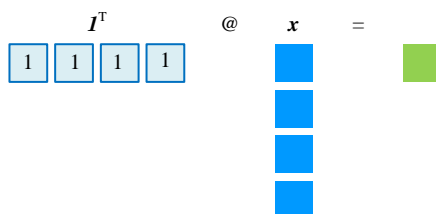


图 4. 求和运算

(6) 除以 n 便是向量 \mathbf{x} 元素平均值：

$$\mathbf{E}(\mathbf{x}) = \frac{x_1 + x_2 + \cdots x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{\mathbf{I} \cdot \mathbf{x}}{n} = \frac{\mathbf{I}^T \mathbf{x}}{n} = \frac{\mathbf{x}^T \mathbf{I}}{n} \quad (7)$$

上式假设前提是， \mathbf{X} 为有 n 个等概率值 $1/n$ 的平均分布。否则，我们要把 $1/n$ 替换成具体的概率值 p_i 。不做特殊说明时，本章默认总体或样本都为等概率。

下式为向量 \mathbf{x} 元素各自平方后再求和：

$$\mathbf{x} \cdot \mathbf{x} = \mathbf{x}^T \mathbf{x} = x_1^2 + x_2^2 + \cdots + x_n^2 = \sum_{i=1}^n x_i^2 \quad (8)$$

上式结果为标量。

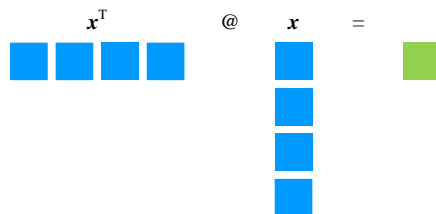


图 5. 平方和运算

计算样本方差时也用到类似 (8) 计算：

$$\text{var}(X) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2 \quad (9)$$

上式中，随机数 X 的样本点构成列向量 \mathbf{x} ， \mathbf{x} 方差则为：

$$\text{var}(\mathbf{x}) = \frac{1}{n-1} \left(\mathbf{x} - \frac{\mathbf{I}^T \mathbf{x}}{n} \right) \cdot \left(\mathbf{x} - \frac{\mathbf{I}^T \mathbf{x}}{n} \right) = \frac{1}{n-1} \left(\mathbf{x} - \frac{\mathbf{I}^T \mathbf{x}}{n} \right)^T \left(\mathbf{x} - \frac{\mathbf{I}^T \mathbf{x}}{n} \right) \quad (10)$$

本书第 22 章将讲解如何展开上式。

前文介绍过，在计算样本协方差时，我们用过类似 (2) 运算：

$$\text{cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - E(X))(y_i - E(Y)) \quad (11)$$

⚠ 注意，如果计算总体方差、协方差的话，(9) 和 (11) 分母的 $n-1$ 则应该改为 n 。当 n 足够大，可以不区分 $n-1$ 或 n 。

上式中，随机数 X 和 Y 的样本点写成列向量 \mathbf{x} 和 \mathbf{y} ，也就是说，(11) 可以写成：

$$\text{cov}(\mathbf{x}, \mathbf{y}) = \frac{1}{n-1} \left(\mathbf{x} - \frac{\mathbf{I}^T \mathbf{x}}{n} \right) \cdot \left(\mathbf{y} - \frac{\mathbf{I}^T \mathbf{y}}{n} \right) = \frac{1}{n-1} \left(\mathbf{x} - \frac{\mathbf{I}^T \mathbf{x}}{n} \right)^T \left(\mathbf{y} - \frac{\mathbf{I}^T \mathbf{y}}{n} \right) \quad (12)$$



统计和线性代数之前有着千丝万缕的联系，本书第 22 章还会继续这一话题。

几何视角

如果 \mathbf{x} 为 n 维单位列向量，则下两式成立：

$$\mathbf{x} \cdot \mathbf{x} = \langle \mathbf{x}, \mathbf{x} \rangle = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|_2^2 = 1, \quad \sqrt{\mathbf{x} \cdot \mathbf{x}} = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\mathbf{x}^T \mathbf{x}} = \|\mathbf{x}\|_2 = 1 \quad (13)$$

整理以上不同等式都得到同一等式：

$$x_1^2 + x_2^2 + \cdots x_n^2 = 1 \quad (14)$$

提醒大家注意，但凡遇到矩阵乘积结果为标量的情况，请考虑是否能从“距离”角度理解这个矩阵乘积。

几何角度，如图 6 (a) 所示，若 $n = 2$ ，(13) 代表平面上的**单位圆** (unit circle)。如图 6 (b) 所示，若 $n = 3$ ，(13) 代表三维空间的**单位球体** (unit sphere)。当 $n > 3$ 时，在多维空间中，(13) 代表 **n 维单位球面** (unit n -sphere) 或**单位超球面** (unit hyper-sphere)。

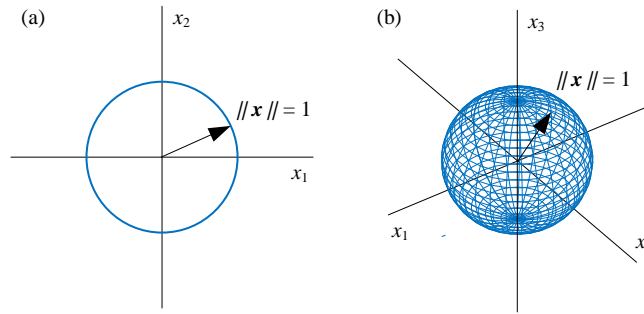


图 6. 单位圆和单位球体

单位圆、单位球、单位超球面内部的点满足：

$$\mathbf{x} \cdot \mathbf{x} = \langle \mathbf{x}, \mathbf{x} \rangle = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|_2^2 < 1, \quad \sqrt{\mathbf{x} \cdot \mathbf{x}} = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\mathbf{x}^T \mathbf{x}} = \|\mathbf{x}\|_2 < 1 \quad (15)$$

即，

$$x_1^2 + x_2^2 + \cdots x_n^2 < 1 \quad (16)$$

单位圆、单位球、单位超球面外部的点满足：

$$\mathbf{x} \cdot \mathbf{x} = \langle \mathbf{x}, \mathbf{x} \rangle = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|_2^2 > 1, \quad \sqrt{\mathbf{x} \cdot \mathbf{x}} = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\mathbf{x}^T \mathbf{x}} = \|\mathbf{x}\|_2 > 1 \quad (17)$$

即，

$$x_1^2 + x_2^2 + \cdots x_n^2 > 1 \quad (18)$$

张量积

列向量 \mathbf{x} 和自身的张量积结果为方阵，相当于 \mathbf{x} 和 \mathbf{x}^T 的乘积：

$$\mathbf{x} \otimes \mathbf{x} = \mathbf{x} @ \mathbf{x}^T = \begin{bmatrix} x_1 x_1 & x_1 x_2 & \cdots & x_1 x_n \\ x_2 x_1 & x_2 x_2 & \cdots & x_2 x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_n x_1 & x_n x_2 & \cdots & x_n x_n \end{bmatrix} \quad (19)$$

图 7 所示为 (19) 计算过程。

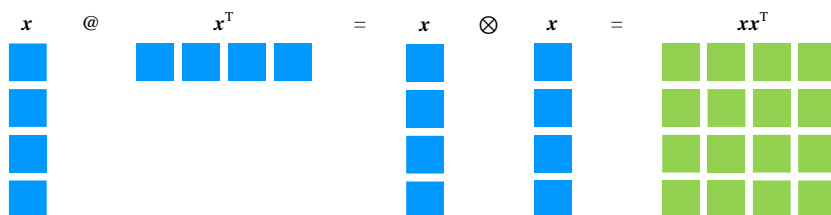


图 7. 张量积运算

用两种方式展开 (19)，可以得到：

$$\begin{aligned} \mathbf{x} \otimes \mathbf{x} = \mathbf{x} \mathbf{x}^T &= \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \mathbf{x}^T = \begin{bmatrix} x_1 \mathbf{x}^T \\ x_2 \mathbf{x}^T \\ \vdots \\ x_n \mathbf{x}^T \end{bmatrix} \\ &= \mathbf{x} \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix} = \begin{bmatrix} x_1 \mathbf{x} & x_2 \mathbf{x} & \cdots & x_n \mathbf{x} \end{bmatrix} \end{aligned} \quad (20)$$

本书前文提过，向量张量积的行向量、列向量都存在“倍数关系”。这实际上解释了为什么非 0 向量张量积的秩 (rank) 为 1。本书第 7 章将介绍“秩”这个概念。另外，请大家注意如图 8 所示的两种形状的张量积和矩阵乘法关系，并注意区分结果形状。

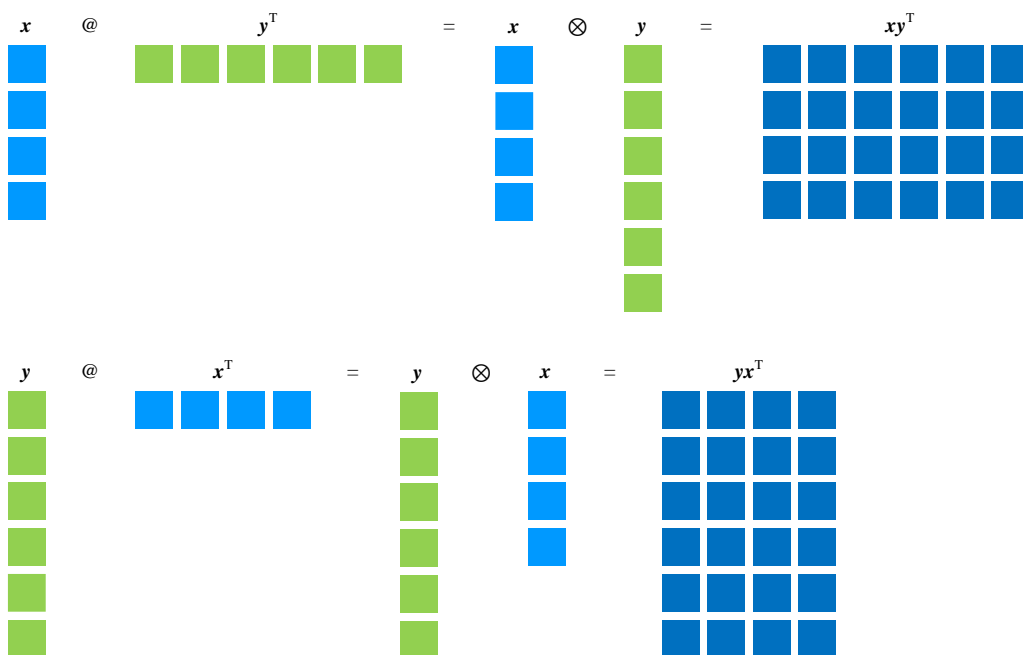


图 8. 另外两种形状的张量积

5.3 再聊全 1 列向量

本节主要介绍全 1 列向量 $\mathbf{1}$ 在求和方面的用途。

有关 Σ 求和，本系列丛书《数学要素》第 14 章中讲过。本节主要从矩阵乘法角度再深入探讨。

每列元素求和

如图 9 所示，全 1 列向量 $\mathbf{1}$ 转置左乘数据矩阵 \mathbf{X} ，相当于对 \mathbf{X} 每一列元素求和。计算结果为行向量，行向量的每个元素是 \mathbf{X} 对应列元素之和：

$$(\mathbf{1}_{n \times 1})^T \mathbf{X} = \begin{bmatrix} 1 & 1 & \cdots & 1 \end{bmatrix}_{1 \times n} \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,D} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,D} \end{bmatrix}_{n \times D} = \begin{bmatrix} \sum_{i=1}^n x_{i,1} & \sum_{i=1}^n x_{i,2} & \cdots & \sum_{i=1}^n x_{i,D} \end{bmatrix}_{1 \times D} \quad (21)$$

请大家格外注意矩阵形状。全 1 列向量 $\mathbf{1}$ 的形状为 $n \times 1$ ，转置之后 $\mathbf{1}^T$ 的形状为 $1 \times n$ 。数据矩阵 \mathbf{X} 的形状为 $n \times D$ 。矩阵乘积 $\mathbf{1}^T \mathbf{X}$ 的结果形状为 $1 \times D$ 。上式就是我们在《数学要素》第 14 章中介绍的“偏求和”的一种。

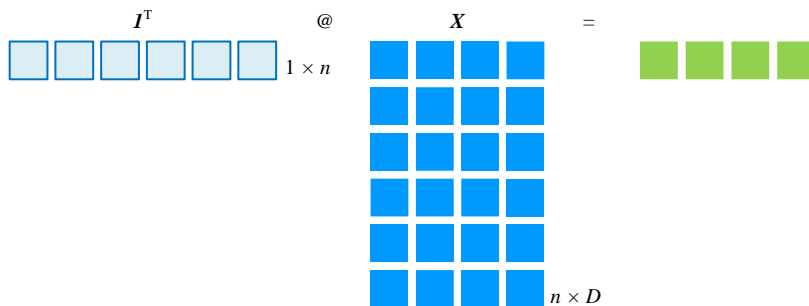


图 9. 列方向求和

(21) 左右除以 n ，便得到每一列元素均值构成的行向量 $E(\mathbf{X})$ ：

$$E(\mathbf{X}) = \frac{\mathbf{1}^T \mathbf{X}}{n} = \begin{bmatrix} \frac{\sum_{i=1}^n x_{i,1}}{n} & \frac{\sum_{i=1}^n x_{i,2}}{n} & \cdots & \frac{\sum_{i=1}^n x_{i,D}}{n} \end{bmatrix} = [\mu_1 \quad \mu_2 \quad \cdots \quad \mu_D] \quad (22)$$

$E(\mathbf{X})$ 常被称作数据矩阵 \mathbf{X} 的**质心** (centroid)。我们也常用 $\boldsymbol{\mu}_X$ 表达质心。 $\boldsymbol{\mu}_X$ 为列向量，是行向量 $E(\mathbf{X})$ 的转置：

$$\boldsymbol{\mu}_X = \mathbf{E}(\mathbf{X})^T = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_D \end{bmatrix} = \frac{\mathbf{X}^T \mathbf{I}}{n} \quad (23)$$

▲ 注意，本系列丛书定义 $\mathbf{E}(\mathbf{X})$ 为行向量。而 $\boldsymbol{\mu}_X$ 为列向量， $\boldsymbol{\mu}_X$ 和 $\mathbf{E}(\mathbf{X})$ 就差在转置上。 $\mathbf{E}(\mathbf{X})$ 一般常配合原始数据矩阵 \mathbf{X} 一起出现，比如利用广播原则去均值。而 $\boldsymbol{\mu}_X$ 多用在分布相关运算中，比如多元高斯分布。

去均值

上一节提到，全 1 列向量有复制的功能。很多应用场合需要将 (22) 复制 n 份，得到一个和原矩阵形状相同的矩阵。下式可以完成这个计算：

$$\mathbf{I}_{n \times 1} @ \mathbf{E}(\mathbf{X})_{1 \times D} = \frac{\mathbf{I}_{n \times 1} \mathbf{I}_{n \times 1}^T \mathbf{X}}{n} = \begin{bmatrix} \mu_1 & \mu_2 & \cdots & \mu_D \\ \mu_1 & \mu_2 & \cdots & \mu_D \\ \vdots & \vdots & \ddots & \vdots \\ \mu_1 & \mu_2 & \cdots & \mu_D \end{bmatrix}_{n \times D} \quad (24)$$

上式结果和数据矩阵 \mathbf{X} 形状一致，都是 $n \times D$ 。其中， $\mathbf{I}_{n \times 1} \mathbf{I}_{n \times 1}^T$ 相当于向量张量积 $\mathbf{I}_{n \times 1} \otimes \mathbf{I}_{n \times 1}$ ，结果为 $n \times n$ 全 1 方阵。利用向量张量积，(24) 可以写成：

$$\mathbf{I}_{n \times 1} @ \mathbf{E}(\mathbf{X})_{1 \times D} = \frac{\mathbf{I}_{n \times 1} \otimes \mathbf{I}_{n \times 1}}{n} \mathbf{X} \quad (25)$$

上式相当于是 \mathbf{X} 向 \mathbf{I} 正交投影，这是本书第 10 章要探讨的内容。

对 \mathbf{X} 去均值 (demean 或 centralize) 就是 \mathbf{X} 的每个元素减去 \mathbf{X} 对应列方向数据均值，即 \mathbf{X} 减去 (24) 得到去均值数据矩阵 \mathbf{X}_c ：

$$\mathbf{X}_c = \mathbf{X} - \frac{\mathbf{I} \mathbf{I}^T \mathbf{X}}{n} = \begin{bmatrix} x_{1,1} - \mu_1 & x_{1,2} - \mu_2 & \cdots & x_{1,D} - \mu_D \\ x_{2,1} - \mu_1 & x_{2,2} - \mu_2 & \cdots & x_{2,D} - \mu_D \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} - \mu_1 & x_{n,2} - \mu_2 & \cdots & x_{n,D} - \mu_D \end{bmatrix}_{n \times D} \quad (26)$$

上式可以整理为：

$$\mathbf{X} - \frac{\mathbf{I} \mathbf{I}^T \mathbf{X}}{n} = \mathbf{I} \mathbf{X} - \frac{\mathbf{I} \mathbf{I}^T \mathbf{X}}{n} = \left(\mathbf{I} - \frac{\mathbf{I} \mathbf{I}^T}{n} \right) \mathbf{X} \quad (27)$$

其中， \mathbf{I} 是单位矩阵，对角线元素都是 1，其余为 0。上式 \mathbf{I} 形状为 $n \times n$ 。



有关去均值运算，本书第 22 章还要深入这一话题。

如图 10 所示，从几何视角来看，去均值相当于将数据的质心平移到原点。为了方便，我们一般利用广播原则计算去均值矩阵 \mathbf{X}_c ，即 $\mathbf{X}_c = \mathbf{X} - \mathbf{E}(\mathbf{X})$ 。

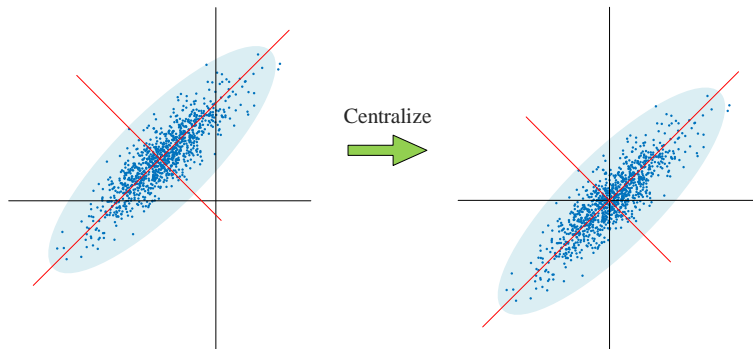


图 10. 去均值的几何视角

用张量积 $\mathbf{I} \otimes \mathbf{I}$ ，(26) 可以写成：

$$\mathbf{X}_c = \mathbf{X} - \frac{\mathbf{I} \otimes \mathbf{I}}{n} \mathbf{X} \quad (28)$$

前文提到，张量积 $\mathbf{I} \otimes \mathbf{I}$ 是个 $n \times n$ 方阵，矩阵的元素都是 1。张量积 $\mathbf{I} \otimes \mathbf{I}$ 再除以 n 得到的方阵每个元素都是 $1/n$ 。

每行元素求和

如图 11 所示，矩阵 \mathbf{X} 乘全 1 列向量 \mathbf{I} ，相当于对 \mathbf{X} 每一行元素求和，结果为列向量：

$$\mathbf{X}_{n \times D} \mathbf{I}_{D \times 1} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,D} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,D} \end{bmatrix}_{n \times D} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{D \times 1} = \begin{bmatrix} \sum_{j=1}^D x_{1,j} \\ \sum_{j=1}^D x_{2,j} \\ \vdots \\ \sum_{j=1}^D x_{n,j} \end{bmatrix}_{n \times 1} \quad (29)$$

⚠ 注意，(22) 和 (29) 两式中的全 \mathbf{I} 向量长度不同。(29) 中全 1 列向量 \mathbf{I} 形状为 $D \times 1$ 。

而 (29) 除以 D 结果是 \mathbf{X} 每行元素平均值，即：

$$\frac{X_{n \times D} \mathbf{I}_{D \times 1}}{D} = \begin{bmatrix} \sum_{j=1}^D x_{1,j} / D \\ \sum_{j=1}^D x_{2,j} / D \\ \vdots \\ \sum_{j=1}^D x_{n,j} / D \end{bmatrix}_{n \times 1} \quad (30)$$

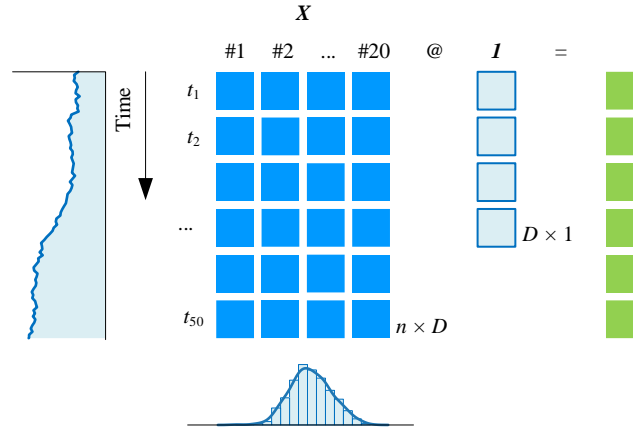


图 11. 行方向求和

大家可能会好奇，数据矩阵的列均值、行均值有怎样的应用场景？

举个例子，假设图 11 中数据矩阵 X 为某个班级 20 名学生一个学期不同时间 t 连续 50 次数学测验成绩。每一列的均值代表的是某个学生的平均成绩，每一行的均值则代表一个班级在某次数学测验的整体表现。采用直方图分析列均值，我们可以得到该学期学生平均成绩的分布。采用线图分析行均值，我们可以得到班级学生平均成绩随时间变化趋势。

所有元素的和

图 12 所示，数据矩阵 X 分别左乘 \mathbf{I}^T 、右乘全 \mathbf{I} 向量，结果为 X 所有元素求和：

$$\mathbf{I}^T X \mathbf{I} = \begin{bmatrix} \sum_{i=1}^n x_{i,1} & \sum_{i=1}^n x_{i,2} & \cdots & \sum_{i=1}^n x_{i,D} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \sum_{j=1}^D \sum_{i=1}^n x_{i,j} \quad (31)$$

上式结果除以 nD ，得到的是整个数据矩阵 X 所有元素的均值。

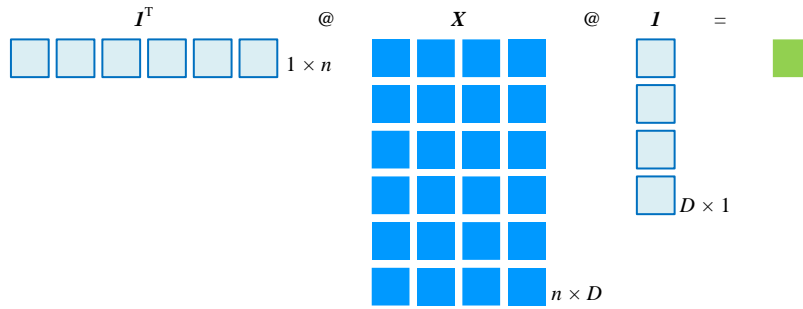


图 12. 矩阵所有元素求和

⚠ 注意，上式中两个全 I 列向量长度也不同，具体形状如图 12 所示。再强调一点，希望大家在看到代数式时，要联想可能的线性代数运算式。本章后续还会继续给出更多示例，以便强化代数和线性代数的联系。

5.4 矩阵乘向量：线性方程组

矩阵 A 为 n 行、 D 列：

$$A_{n \times D} = \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,D} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1} & a_{n,2} & \cdots & a_{n,D} \end{bmatrix} \quad (32)$$

\mathbf{x} 为 D 个未知量 $x_1, x_1 \dots x_D$ 构成的列向量， \mathbf{b} 为 n 个常数 $b_1, b_1 \dots b_n$ 构成的列向量：

$$\mathbf{x}_{D \times 1} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}, \quad \mathbf{b}_{n \times 1} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \quad (33)$$

如图 13 所示， $A\mathbf{x} = \mathbf{b}$ 可以写成：

$$\underbrace{\begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,D} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1} & a_{n,2} & \cdots & a_{n,D} \end{bmatrix}}_{A_{n \times D}} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}}_{\mathbf{x}_{D \times 1}} = \underbrace{\begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}}_{\mathbf{b}_{n \times 1}} \quad (34)$$

(34) 展开得到线性方程组 (system of linear equations)：

$$\begin{cases} a_{1,1}x_1 + a_{1,2}x_2 + \cdots + a_{1,D}x_D = b_1 \\ a_{2,1}x_1 + a_{2,2}x_2 + \cdots + a_{2,D}x_D = b_2 \\ \vdots \\ a_{n,1}x_1 + a_{n,2}x_2 + \cdots + a_{n,D}x_D = b_n \end{cases} \quad (35)$$

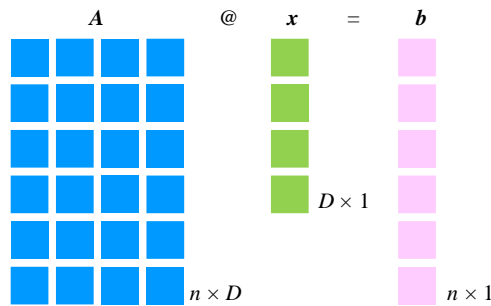


图 13. 长方阵乘列向量

解的个数

若 (34) 有唯一一组解，矩阵 A 可逆，即：

$$Ax = b \Rightarrow x = A^{-1}b \quad (36)$$

此时称 $Ax = b$ 为恰定方程组。

有无穷多解的方程组被称作**欠定方程组** (underdetermined system)。

解不存在的方程组被称作**超定方程组** (overdetermined system)。

特别地，如果 $A^T A$ 可逆， x 可以通过下式求解：

$$Ax = b \Rightarrow A^T A x = A^T b \Rightarrow x = \underbrace{(A^T A)^{-1}}_{A^+} A^T b \quad (37)$$

$(A^T A)^{-1} A^T$ 常被称作**广义逆** (generalized inverse)，或**伪逆** (pseudoinverse)。

➡ 本系列丛书《数学要素》一册介绍过**最小二乘法** (ordinary least squares, OLS) 和广义逆之间的关系。本系列丛书《概率统计》和《数据科学》两册还会深入讲解最小二乘法回归。

线性代数本身具有“代数”属性，这也就是为什么很多教材以求解 $Ax = b$ 为起点讲解线性代数。而本书则试图跳出“代数”的桎梏，从向量、几何、空间、数据等视角理解 $Ax = b$ 。

线性组合视角

下面用另外一个视角看 $Ax = b$ 。

本书前文反复提到，矩阵 A 可以看做由一组列向量构造而成：

$$A_{n \times D} = [a_1 \quad a_2 \quad \cdots \quad a_D] \quad (38)$$

如图 14 所示, (34) 可以写成:

$$\begin{bmatrix} \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_D \end{bmatrix}_{1 \times D} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}_{D \times 1} = \mathbf{b}_{n \times 1} \quad (39)$$

$$\begin{array}{c} \mathbf{A} \\ \mathbf{a}_1 \quad \mathbf{a}_2 \quad \cdots \quad \mathbf{a}_D \end{array} \quad 1 \times D \quad @ \quad \begin{array}{c} \mathbf{x} \\ \text{ } \\ \text{ } \\ \text{ } \\ \text{ } \end{array} \quad D \times 1 = \mathbf{b}$$

图 14. 线性组合视角看线性方程组

展开 (39) 得到:

$$x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + \cdots + x_D \mathbf{a}_D = \mathbf{b}_{n \times 1} \quad (40)$$

即,

$$x_1 \begin{bmatrix} a_{1,1} \\ a_{2,1} \\ \vdots \\ a_{n,1} \end{bmatrix} + x_2 \begin{bmatrix} a_{1,2} \\ a_{2,2} \\ \vdots \\ a_{n,2} \end{bmatrix} + \cdots + x_D \begin{bmatrix} a_{1,D} \\ a_{2,D} \\ \vdots \\ a_{n,D} \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \quad (41)$$

$\mathbf{a}_1 \quad \mathbf{a}_2 \quad \mathbf{a}_D$

当 x_1, x_2, \dots, x_D 取具体值时, 上式代表**线性组合** (linear combination)。用腊八粥举个例子, 上式相当于不同比例的原料混合, x_i 就是比例, \mathbf{a}_i 就是不同的原料。而 \mathbf{b} 就是混合得到的八宝粥。



线性组合这个概念非常重要, 本书第 7 章将专门介绍。

映射视角

如图 15 所示, 从**线性映射** (linear mapping) 角度来看, (34) 代表从 \mathbb{R}^D 空间到 \mathbb{R}^n 空间的某种特定映射。列向量 $\mathbf{x}_{D \times 1}$ 在 \mathbb{R}^D 中, 而列向量 $\mathbf{b}_{n \times 1}$ 在 \mathbb{R}^n 中。当且仅当矩阵 \mathbf{A} 可逆, 可以完成从 \mathbb{R}^D 空间到 \mathbb{R}^n 空间的映射。这种情况下, $n = D$, 也就是两个空间相同, 我们管这种线性映射叫**线性变换** (linear transformation)。

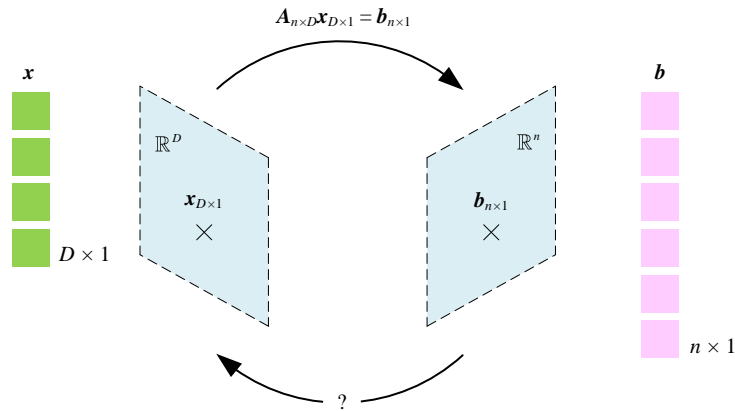


图 15. 线性映射

几何视角

如果 2 维向量 $\mathbf{x} = [x_1, x_2]^T$ 的模为 1， \mathbf{x} 的起点位于原点，终点则位于单位圆上。给定如下矩阵 S 和 R ：

$$S = \begin{bmatrix} 2 & \\ & 1 \end{bmatrix}, \quad R = \begin{bmatrix} \sqrt{2}/2 & -\sqrt{2}/2 \\ \sqrt{2}/2 & \sqrt{2}/2 \end{bmatrix} \quad (42)$$

利用矩阵乘法， \mathbf{x} 分别经过 S 和 R ($A = RS$) 映射得到 \mathbf{y} ：

$$\mathbf{y} = \mathbf{Ax} = \mathbf{RSx} = \underbrace{\begin{bmatrix} \sqrt{2}/2 & -\sqrt{2}/2 \\ \sqrt{2}/2 & \sqrt{2}/2 \end{bmatrix}}_R \underbrace{\begin{bmatrix} 2 \\ 1 \end{bmatrix}}_S \mathbf{x} \quad (43)$$

如图 16 所示，(43) 代表“缩放 → 旋转”。请大家注意几何变换的先后顺序，缩放 (S) 先作用于 \mathbf{x} ，对应矩阵乘法 $S\mathbf{x}$ ；然后，旋转 (R) 再作用于 $S\mathbf{x}$ ，得到 $RS\mathbf{x}$ 。准确来说，图 16 中的这两种几何变换叫做线性变换，这是本书第 8 章要探讨的话题。

也就是说，矩阵连乘代表一系列有先后顺序的几何变换。此外，以上分析还告诉我们矩阵 A 可以分解为 S 和 R 相乘，用到的数学工具就是第 11 章要讲的矩阵分解。

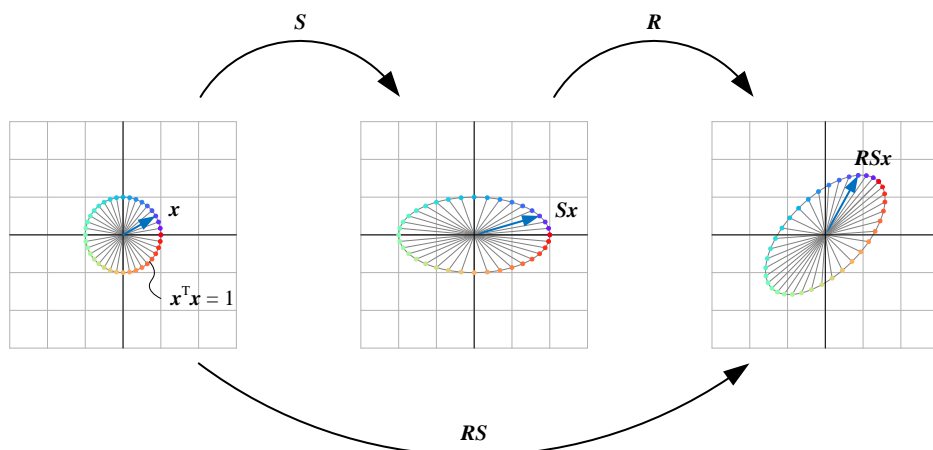


图 16. 几何变换视角

向量模

凡是向量就有自己的长度，即向量模、 L^2 范数。 $b_{n \times 1}$ 的向量模、 L^2 范数为：

$$\|b\| = \|Ax\| \quad (44)$$

注意， b 的模是标量。

利用矩阵乘法，上式可以写成：

$$\|b\| = \sqrt{b^T b} = \sqrt{x^T A^T A x} \quad (45)$$

b 的模的平方则为：

$$\|b\|_2^2 = b^T b = x^T A^T A x \quad (46)$$

$x^T A^T A x$ 这种矩阵乘法的结果为非负标量，其中 $A^T A$ 叫做 A 的格拉姆矩阵。 $x^T A^T A x$ 就是下一节要介绍的二次型。

举个例子，如果向量 x 的模为1，平面上向量 x 的终点在单位圆上。如图17所示，经过如下 $Ax = b$ 的线性映射得到的向量 b 终点在旋转椭圆上：

$$b = Ax = \underbrace{\begin{bmatrix} 1.25 & -0.75 \\ -0.75 & 1.25 \end{bmatrix}}_A x \quad (47)$$

而矩阵 A 恰好可逆，通过如下运算，我们把旋转椭圆变换成单位圆：

$$x = A^{-1}b = \underbrace{\begin{bmatrix} 1.25 & 0.75 \\ 0.75 & 1.25 \end{bmatrix}}_{A^{-1}} b \quad (48)$$

大家可能会好奇，该如何计算旋转椭圆的半长轴、半短轴的长度，以及长轴旋转角度等信息？本书第 14 章将给出答案。

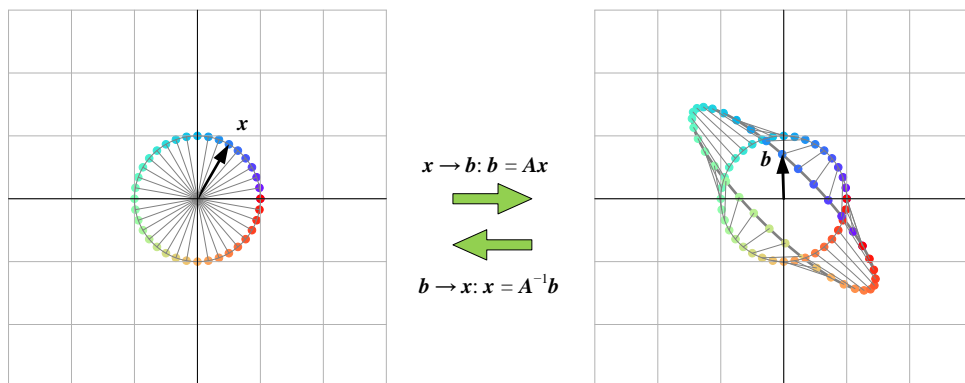


图 17. 单位圆到旋转椭圆

5.5 向量乘矩阵乘向量：二次型

二次型 (quadratic form) 的矩阵算式为：

$$\mathbf{x}^T \mathbf{Q} \mathbf{x} = q \quad (49)$$

其中， \mathbf{Q} 为对称阵， q 为实数。 \mathbf{Q} 和 \mathbf{x} 分别为：

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}, \quad \mathbf{Q} = \begin{bmatrix} q_{1,1} & q_{1,2} & \cdots & q_{1,D} \\ q_{2,1} & q_{2,2} & \cdots & q_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ q_{D,1} & q_{D,2} & \cdots & q_{D,D} \end{bmatrix} \quad (50)$$

(49) 对应的矩阵运算过程如图 18 所示。

$\mathbf{x}^T \mathbf{Q} \mathbf{x}$ 像极了 $\mathbf{x}^T \mathbf{x}$ ，也就是说 $\mathbf{x}^T \mathbf{Q} \mathbf{x}$ 类似 $\|\mathbf{x}\|_2^2$ ，结果都是“标量”。几何角度， $\|\mathbf{x}\|_2^2$ 代表向量 \mathbf{x} 度的平方， $\mathbf{x}^T \mathbf{Q} \mathbf{x}$ 似乎也代表着某种“距离的平方”，本书后续将会专门介绍。

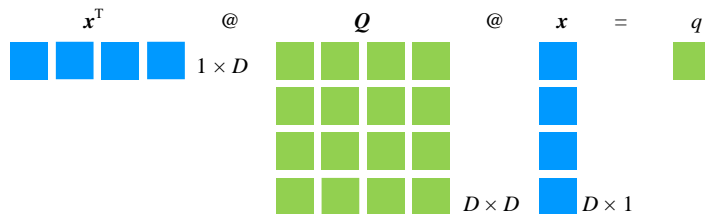


图 18. $\mathbf{x}^T \mathbf{Q} \mathbf{x} = q$ 矩阵运算

将 (50) 代入 (49)，展开得到：

$$\mathbf{x}^T \mathbf{Q} \mathbf{x} = \sum_{i=1}^D q_{i,i} x_i^2 + \sum_{i=1}^D \sum_{j=1}^D q_{i,j} x_i x_j = q, \quad i \neq j \quad (51)$$

观察上式，发现单项式变量的最高次数为 2，这就是为什么 $\mathbf{x}^T \mathbf{Q} \mathbf{x}$ 叫二次型的原因。

举个例子

比如 \mathbf{x} 和 \mathbf{Q} 分别为：

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \mathbf{Q} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad (52)$$

代入 (49) 得到：

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = ax_1^2 + (b+c)x_1x_2 + dx_2^2 = q \quad (53)$$

可以发现，(53) 对应本系列丛书中《数学要素》介绍过各种二次曲线，比如正圆、椭圆、抛物线或双曲线，具体如图 19 所示。



本书第 20 章还要用线性代数工具深入探讨这些圆锥曲线。

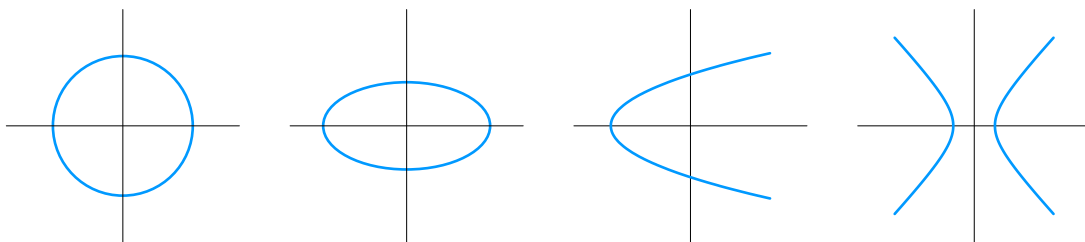


图 19. 四种二次曲线

将 (53) 写成二元函数形式 $f(x_1, x_2)$ ：

$$f(x_1, x_2) = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = ax_1^2 + (b+c)x_1x_2 + dx_2^2 \quad (54)$$

(54) 对应着如图 20 所示的几种曲面。而 $f(x_1, x_2) = q$ ，相当于曲面某个高度的等高线。

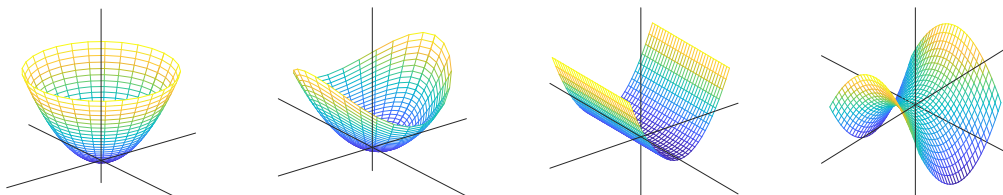


图 20. 常见二次型曲面



本书第 21 章将探讨图 20 这些曲面和正定性、极值之间的联系。

高斯分布

二次型的应用无处不在。举个例子，二元正态分布的概率密度函数解析式如下：

$$f_{X_1, X_2}(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho_{1,2}^2}} \times \exp \left(\frac{-1}{2} \left(\frac{1}{(1-\rho_{1,2}^2)} \left(\left(\frac{x_1-\mu_1}{\sigma_1} \right)^2 - 2\rho_{1,2} \left(\frac{x_1-\mu_1}{\sigma_1} \right) \left(\frac{x_2-\mu_2}{\sigma_2} \right) + \left(\frac{x_2-\mu_2}{\sigma_2} \right)^2 \right) \right) \right) \quad (55)$$

大家应该记得我们在本系列丛书《数学要素》第 9 章介绍过这种形式椭圆。

而多元正态分布的概率密度函数为：

$$f_{\mathbf{x}}(\mathbf{x}) = \frac{\exp \left(-\frac{1}{2} \overbrace{\left(\mathbf{x} - \boldsymbol{\mu} \right)^T \boldsymbol{\Sigma}^{-1} \left(\mathbf{x} - \boldsymbol{\mu} \right)}^{\text{Ellipse}} \right)}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \quad (56)$$

分子中已经明显看到类似 (49) 的矩阵乘法。本书第 20 章会继续这一话题。

比较上两式，大家也应该清楚，为什么进入多元领域，比如多元微积分、多元概率统计，我们便离不开线性代数。二元正态分布的概率密度函数解析式已经如此复杂，更不用说三元、四元，乃至 D 元。

三个方阵连乘

我们再看另外矩阵乘法一种形式，具体如下：

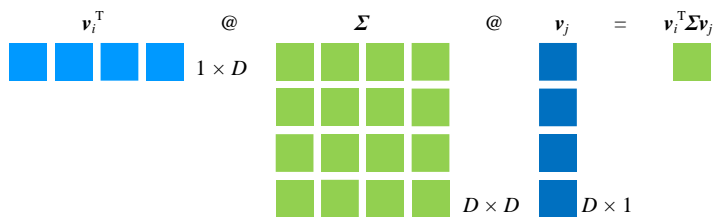
$$\mathbf{V}^T \boldsymbol{\Sigma} \mathbf{V} \quad (57)$$

其中， \mathbf{V} 和 $\boldsymbol{\Sigma}$ 都是 $D \times D$ 方阵，上式结果也是 $D \times D$ 方阵。特别地，实际应用中 \mathbf{V} 多为正交矩阵，即为 \mathbf{V} 方阵且满足 $\mathbf{V}\mathbf{V}^T = \mathbf{I}$ 。

将 \mathbf{V} 写成 $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_D]$ ，展开 (57) 得到：

$$\begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_D^T \end{bmatrix} \boldsymbol{\Sigma} \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_D \end{bmatrix} = \begin{bmatrix} \mathbf{v}_1^T \boldsymbol{\Sigma} \mathbf{v}_1 & \mathbf{v}_1^T \boldsymbol{\Sigma} \mathbf{v}_2 & \cdots & \mathbf{v}_1^T \boldsymbol{\Sigma} \mathbf{v}_D \\ \mathbf{v}_2^T \boldsymbol{\Sigma} \mathbf{v}_1 & \mathbf{v}_2^T \boldsymbol{\Sigma} \mathbf{v}_2 & \cdots & \mathbf{v}_2^T \boldsymbol{\Sigma} \mathbf{v}_D \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{v}_D^T \boldsymbol{\Sigma} \mathbf{v}_1 & \mathbf{v}_D^T \boldsymbol{\Sigma} \mathbf{v}_2 & \cdots & \mathbf{v}_D^T \boldsymbol{\Sigma} \mathbf{v}_D \end{bmatrix} \quad (58)$$

结果中，矩阵 (i, j) 元素 $\mathbf{v}_i^T \boldsymbol{\Sigma} \mathbf{v}_j$ 便是一个二次型， $\mathbf{v}_i^T \boldsymbol{\Sigma} \mathbf{v}_j$ 对应的运算示意图如图 21 所示。这说明，上式包含了 $D \times D$ 个二次型。

图 21. $v_i^T \Sigma v_j$ 矩阵运算

➡ 二次型在多元微积分、正定性、多元正态分布、协方差矩阵、数据映射和优化方法中都有举足轻重的分量。本书后续将会深入探讨。

5.6 方阵乘方阵：矩阵分解

和方阵有关的矩阵乘法中，方阵乘方阵最为简单。图 22 所示两种方阵乘法常见于 LU 分解、Cholesky 分解、特征值分解等场合。

➡ 本节不展开讲解矩阵分解，本书第 11 ~ 16 章将专门介绍不同类别矩阵分解。

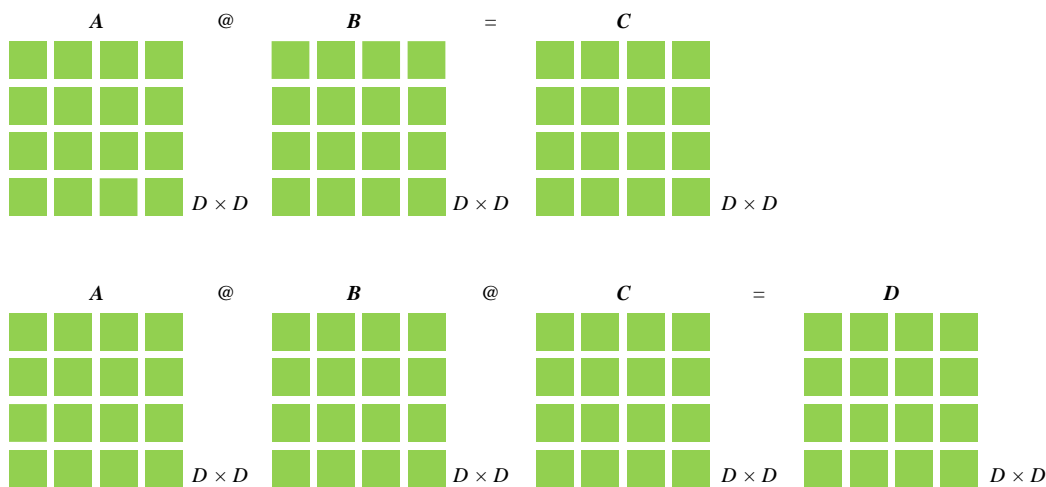


图 22. 方阵乘方阵

特别地，方阵 A 如果满足：

$$A^2 = A \quad (59)$$

则称 A 为**幂等矩阵** (idempotent matrix)。



我们会在本书统计部分和最小二乘法线性回归中再谈及幂等矩阵。此外，丛书每册均有涉及线性回归这个话题，本书采用的是线性代数和向量几何视角，《概率统计》则利用统计视角理解线性回归，而《数据科学》则是从数据分析视角介绍如何应用这个模型。

5.7 对角阵：批量缩放

如果形状相同的方阵 \mathbf{A} 和 \mathbf{B} 都为对角阵，两者乘积还是一个对角阵：

$$\mathbf{A}_{D \times D} \mathbf{B}_{D \times D} = \begin{bmatrix} a_1 & & & \\ & a_2 & & \\ & & \ddots & \\ & & & a_D \end{bmatrix} \begin{bmatrix} b_1 & & & \\ & b_2 & & \\ & & \ddots & \\ & & & b_D \end{bmatrix} = \begin{bmatrix} a_1 b_1 & & & \\ & a_2 b_2 & & \\ & & \ddots & \\ & & & a_D b_D \end{bmatrix} \quad (60)$$

对角阵 \mathbf{A} 的逆也是一个对角阵：

$$\mathbf{A}_{D \times D} (\mathbf{A}_{D \times D})^{-1} = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_D \end{bmatrix} \begin{bmatrix} 1/\lambda_1 & & & \\ & 1/\lambda_2 & & \\ & & \ddots & \\ & & & 1/\lambda_D \end{bmatrix} = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} = \mathbf{I}_{D \times D} \quad (61)$$

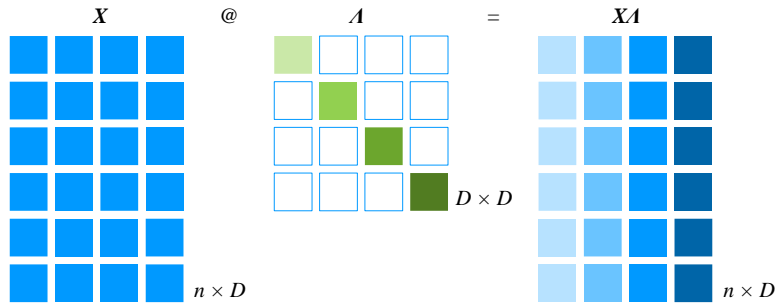
其中， $\lambda_j \neq 0$ 。注意，本书中经常采用 $\mathbf{\Lambda}$ (capital lambda) 和 \mathbf{S} 代表对角阵。

右乘

矩阵 \mathbf{X} 乘 $D \times D$ 对角方阵 \mathbf{A} ：

$$\begin{aligned} \mathbf{X}_{n \times D} \mathbf{A}_{D \times D} &= \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_D \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_D \end{bmatrix} \\ &= \begin{bmatrix} \lambda_1 \mathbf{x}_1 & \lambda_2 \mathbf{x}_2 & \cdots & \lambda_D \mathbf{x}_D \end{bmatrix} \end{aligned} \quad (62)$$

观察 (62) 发现， \mathbf{A} 的对角线元素相当于缩放系数，分别对矩阵 \mathbf{X} 的每一列数值进行不同比例缩放。

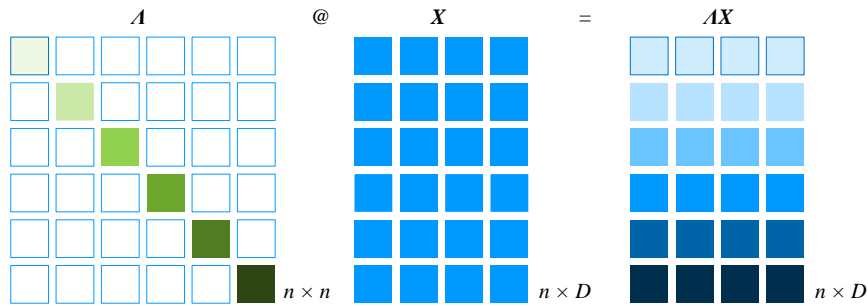
图 23. X 乘对角方阵 A

左乘

$n \times n$ 对角阵 A 左乘矩阵 X :

$$A_{n \times n} X_{n \times D} = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix}_{n \times n} \begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \\ \vdots \\ \mathbf{x}^{(n)} \end{bmatrix}_{n \times 1} = \begin{bmatrix} \lambda_1 \mathbf{x}^{(1)} \\ \lambda_2 \mathbf{x}^{(2)} \\ \vdots \\ \lambda_n \mathbf{x}^{(n)} \end{bmatrix}_{n \times 1} \quad (63)$$

观察 (63)，可以发现 A 的对角线元素分别对矩阵 X 的每一行数值进行批量缩放。

图 24. 对角阵 A 乘方阵 X

乘行向量

特别地，行向量 $\mathbf{x}^{(i)}$ 乘 $D \times D$ 对角阵 A ，相当于对行向量每个元素以不同比例分别缩放：

$$\mathbf{x}^{(i)} A_{D \times D} = \begin{bmatrix} x_{1,i} & x_{1,2} & \cdots & x_{1,D} \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_D \end{bmatrix}_{D \times D} = \begin{bmatrix} \lambda_1 x_{1,i} & \lambda_2 x_{1,2} & \cdots & \lambda_D x_{1,D} \end{bmatrix} \quad (64)$$

乘列向量

类似地， $n \times n$ 对角阵 A 乘列向量 \mathbf{x} ，相当于对列向量每个元素以不同比例分别缩放：

$$A_{n \times n} \mathbf{x}_{n \times 1} = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix}_{n \times n} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} \lambda_1 x_1 \\ \lambda_2 x_2 \\ \vdots \\ \lambda_n x_n \end{bmatrix}_{n \times 1} \quad (65)$$

左右都乘

再看下例， $D \times D$ 对角方阵 A 分别左乘、右乘 $D \times D$ 方阵 B ：

$$\begin{aligned} ABA &= \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_D \end{bmatrix} \begin{bmatrix} b_{1,1} & b_{1,2} & \cdots & b_{1,D} \\ b_{2,1} & b_{2,2} & \cdots & b_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ b_{D,1} & b_{D,2} & \cdots & b_{D,D} \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_D \end{bmatrix} \\ &= \begin{bmatrix} \lambda_1 \lambda_1 b_{1,1} & \lambda_1 \lambda_2 b_{1,2} & \cdots & \lambda_1 \lambda_D b_{1,D} \\ \lambda_2 \lambda_1 b_{2,1} & \lambda_2 \lambda_2 b_{2,2} & \cdots & \lambda_2 \lambda_D b_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_D \lambda_1 b_{D,1} & \lambda_D \lambda_2 b_{D,2} & \cdots & \lambda_D \lambda_D b_{D,D} \end{bmatrix} \end{aligned} \quad (66)$$

看到 (66) 结果形式，大家是否想到了协方差矩阵。 λ_i 相当于均方差， $b_{i,j}$ 相当于相关性系数。

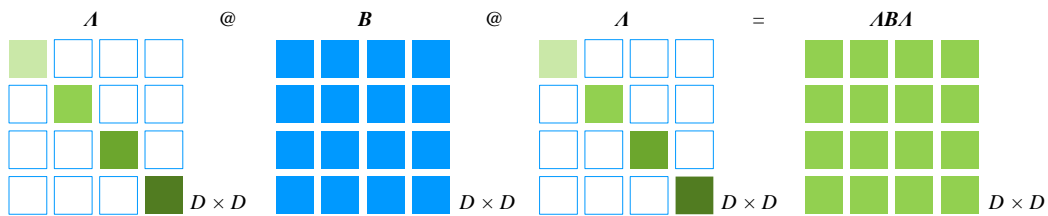


图 25. 对角阵 A 分别左乘、右乘方阵 B

二次型特例

再看一个二次型的特例：

$$\mathbf{x}^T A_{D \times D} \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}^T \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_D \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \lambda_1 x_1^2 + \lambda_2 x_2^2 + \cdots + \lambda_D x_D^2 = \sum_{j=1}^D \lambda_j x_j^2 \quad (67)$$

图 26 所示为上述运算的示意图。

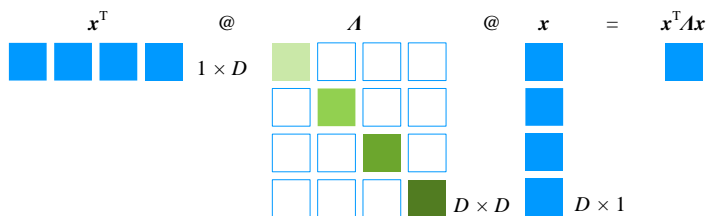


图 26. $\mathbf{x}^T \mathbf{A} \mathbf{x}$ 对应的矩阵运算

几何视角

看到类似 (67) 形式运算，希望大家能联想到正椭圆、正椭球、正椭圆抛物面。比如，如果 $\lambda_1 > \lambda_2 > 0$ ，且 $k > 0$ ，下式对应正椭圆：

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = k \quad (68)$$

这个椭圆的半长轴长度为 $\sqrt{k/\lambda_2}$ ，半短轴长度为 $\sqrt{k/\lambda_1}$ 。

举个例子，下式对应的正椭圆半长轴长度为 2，半短轴长度为 1：

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 1/4 & \\ & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \frac{1}{4} x_1^2 + x_2^2 = 1 \quad (69)$$

再次强调，如果在矩阵运算时遇到对角阵，请试着从几何体缩放角度来看。

5.8 置换矩阵：调换元素顺序

行向量 \mathbf{a} 乘副对角矩阵，如果副对角线上元素都为 1，得到左右翻转的行向量：

$$\begin{bmatrix} a_1 & a_2 & \cdots & a_D \end{bmatrix}_{1 \times D} \begin{bmatrix} & & & 1 \\ & & 1 & \\ & \ddots & & \\ 1 & & & \end{bmatrix}_{D \times D} = \begin{bmatrix} a_D & a_{D-1} & \cdots & a_1 \end{bmatrix} \quad (70)$$

实际上，(70) 中完成左右翻转的方阵是**置换矩阵** (permutation matrix) 的一种特殊形式。

置换矩阵是由 0 和 1 组成的方阵。置换矩阵的每一行、每一列都恰好只有一个 1，其余元素均为 0。置换矩阵的作用是调换元素顺序。

举个例子：

$$\begin{bmatrix} a_1 & a_2 & a_3 & a_4 \end{bmatrix} \begin{bmatrix} & & 1 & \\ & & & 1 \\ 1 & & & \\ & & & & 1 \end{bmatrix} = \begin{bmatrix} a_3 & a_1 & a_4 & a_2 \end{bmatrix} \quad (71)$$

调整列向量顺序

置换矩阵同样可以作用于矩阵，将 (71) 中行向量元素替换成列向量，即，

$$\mathbf{a}_1 = \begin{bmatrix} a_{1,1} \\ a_{2,1} \\ a_{3,1} \\ a_{4,1} \end{bmatrix}, \quad \mathbf{a}_2 = \begin{bmatrix} a_{1,2} \\ a_{2,2} \\ a_{3,2} \\ a_{4,2} \end{bmatrix}, \quad \mathbf{a}_3 = \begin{bmatrix} a_{1,3} \\ a_{2,3} \\ a_{3,3} \\ a_{4,3} \end{bmatrix}, \quad \mathbf{a}_4 = \begin{bmatrix} a_{1,4} \\ a_{2,4} \\ a_{3,4} \\ a_{4,4} \end{bmatrix} \quad (72)$$

可以得到：

$$\begin{bmatrix} \begin{bmatrix} a_{1,1} \\ a_{2,1} \\ a_{3,1} \\ a_{4,1} \end{bmatrix} & \begin{bmatrix} a_{1,2} \\ a_{2,2} \\ a_{3,2} \\ a_{4,2} \end{bmatrix} & \begin{bmatrix} a_{1,3} \\ a_{2,3} \\ a_{3,3} \\ a_{4,3} \end{bmatrix} & \begin{bmatrix} a_{1,4} \\ a_{2,4} \\ a_{3,4} \\ a_{4,4} \end{bmatrix} \end{bmatrix} \begin{bmatrix} & & 1 & \\ & & & 1 \\ 1 & & & \\ & & & & 1 \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} a_{1,3} \\ a_{2,3} \\ a_{3,3} \\ a_{4,3} \end{bmatrix} & \begin{bmatrix} a_{1,1} \\ a_{2,1} \\ a_{3,1} \\ a_{4,1} \end{bmatrix} & \begin{bmatrix} a_{1,4} \\ a_{2,4} \\ a_{3,4} \\ a_{4,4} \end{bmatrix} & \begin{bmatrix} a_{1,2} \\ a_{2,2} \\ a_{3,2} \\ a_{4,2} \end{bmatrix} \end{bmatrix} \quad (73)$$

大家看到置换矩阵右乘矩阵 \mathbf{A} ，让 \mathbf{A} 的列向量顺序改变。

调整行向量顺序

这个置换矩阵左乘矩阵 \mathbf{A} ，可以改变 \mathbf{A} 的行向量的排序：

$$\begin{bmatrix} & & 1 & \\ & & & 1 \\ 1 & & & \\ & & & & 1 \end{bmatrix} \begin{bmatrix} \mathbf{a}^{(1)} \\ \mathbf{a}^{(2)} \\ \mathbf{a}^{(3)} \\ \mathbf{a}^{(4)} \end{bmatrix} = \begin{bmatrix} \mathbf{a}^{(2)} \\ \mathbf{a}^{(4)} \\ \mathbf{a}^{(1)} \\ \mathbf{a}^{(3)} \end{bmatrix} \quad (74)$$

置换矩阵可以用来简化一些矩阵运算。

5.9 矩阵乘向量：映射到一维

前文提到过，任何矩阵乘法都可以从线性映射 (linear mapping) 角度理解。本节和下一节专门从几何角度聊聊线性映射。

形状为 $n \times D$ 矩阵 \mathbf{X} 乘 $D \times 1$ 列向量 \mathbf{v} 得到 $n \times 1$ 列向量 \mathbf{z} ：

$$\mathbf{X}_{n \times D} \mathbf{v}_{D \times 1} = \mathbf{z}_{n \times 1} \quad (75)$$

如图 27 所示，矩阵 \mathbf{X} 有 D 列，对应 D 个特征。而结果 \mathbf{z} 只有一列，也就是一个特征。类似 (41)，(75) 也可以写成“线性组合”：

$$\underbrace{\begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_D \end{bmatrix}}_{\mathbf{X}} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_D \end{bmatrix}_{\mathbf{v}} = v_1 \mathbf{x}_1 + v_2 \mathbf{x}_2 + \cdots + v_D \mathbf{x}_D = \mathbf{z} \quad (76)$$

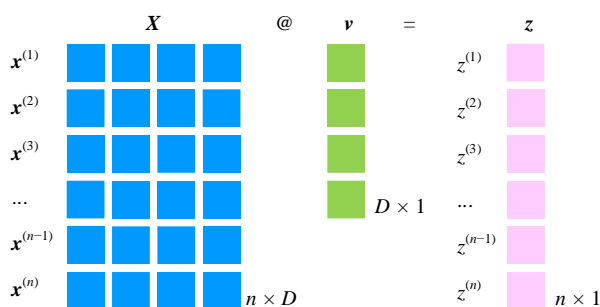


图 27. 矩阵乘法 $\mathbf{X}\mathbf{v} = \mathbf{z}$

此外， $\mathbf{X}\mathbf{v} = \mathbf{z}$ 可以展开写成：

$$\mathbf{X}\mathbf{v} = \begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \\ \vdots \\ \mathbf{x}^{(n)} \end{bmatrix} \mathbf{v} = \begin{bmatrix} \mathbf{x}^{(1)}\mathbf{v} \\ \mathbf{x}^{(2)}\mathbf{v} \\ \vdots \\ \mathbf{x}^{(n)}\mathbf{v} \end{bmatrix} = \begin{bmatrix} z^{(1)} \\ z^{(2)} \\ \vdots \\ z^{(n)} \end{bmatrix} \quad (77)$$

几何视角来看，矩阵 \mathbf{X} 中任意一行 $\mathbf{x}^{(i)}$ 看做是多维坐标系的一个点，运算 $\mathbf{x}^{(i)}\mathbf{v}$ 则是点 $\mathbf{x}^{(i)}$ 在 \mathbf{v} 方向映射， $z^{(i)}$ 则是结果在 \mathbf{v} 上的坐标。如图 28 所示，(75) 这个矩阵乘法运算过程相当于降维。

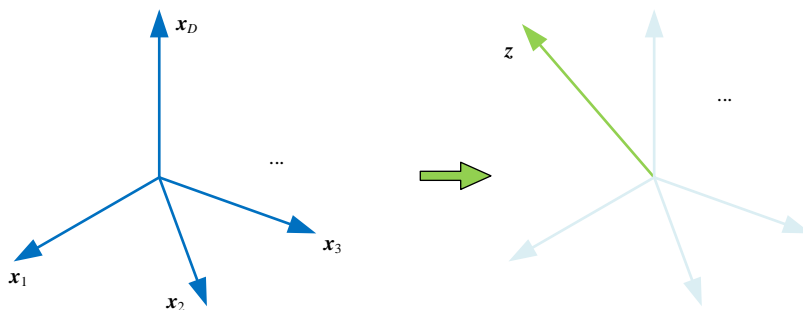


图 28. 多维到一维映射

以鸢尾花数据为例

为了方便理解，下面我们将给 \mathbf{v} 赋予具体数值来讲解。

以鸢尾花数据为例，矩阵 \mathbf{X} 的 4 列分别对应 4 个特征——花萼长度、花萼宽度、花瓣长度、花瓣宽度。 $\mathbf{X}\mathbf{v}=\mathbf{z}$ 结果只有 1 列，相当于只有 1 个特征。

举个例子，如果 \mathbf{v} 中第三个元素为 1，其余元素均为 0，如图 29 所示。向量乘积 $\mathbf{X}\mathbf{v}$ 的结果是从 \mathbf{X} 中提取第 3 列 \mathbf{x}_3 ：

$$\mathbf{X}\mathbf{v} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \mathbf{x}_3 & \mathbf{x}_4 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} = \mathbf{x}_3 \quad (78)$$

也就是说，运算结果只保留第三列花瓣长度相关数据。

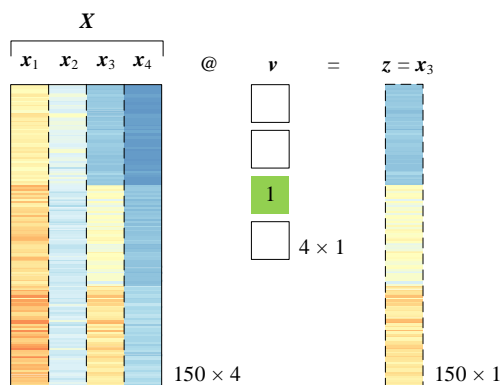


图 29. \mathbf{v} 只有第三个元素为 1，其余均为 0

再举个例子，若我们想要计算每个样本花萼长度 (\mathbf{x}_1)、花萼宽度 (\mathbf{x}_2) 的平均值，可以通过如下运算得到：

$$\mathbf{X}\mathbf{v} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \mathbf{x}_3 & \mathbf{x}_4 \end{bmatrix} \begin{bmatrix} 1/2 \\ 1/2 \\ 0 \\ 0 \end{bmatrix} = \frac{\mathbf{x}_1 + \mathbf{x}_2}{2} \quad (79)$$

同理，下式计算每个样本花萼长度 (\mathbf{x}_1)、花萼宽度 (\mathbf{x}_2)、花瓣长度 (\mathbf{x}_3)、花瓣宽度 (\mathbf{x}_4) 四个特征平均值：

$$\mathbf{X}\mathbf{v} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \mathbf{x}_3 & \mathbf{x}_4 \end{bmatrix} \begin{bmatrix} 1/4 \\ 1/4 \\ 1/4 \\ 1/4 \end{bmatrix} = \frac{\mathbf{x}_1 + \mathbf{x}_2 + \mathbf{x}_3 + \mathbf{x}_4}{4} \quad (80)$$

几何角度来看上式运算，(80) 相当于 4 维空间的散点，被压缩到了一条轴上，具体如图 30。图 30 中 4 维空间的散点仅仅是示意图而已。

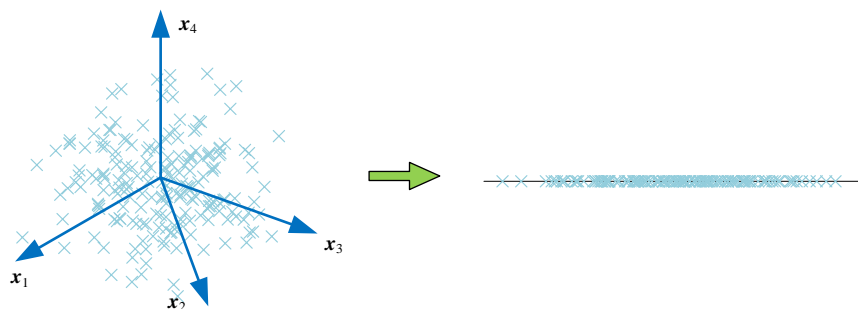


图 30. 4 维空间散点压缩到 1 维

5.10 矩阵乘矩阵：映射到多维

有了上一节内容做基础，这一节我们介绍矩阵乘法在多维映射中扮演色的角色。

两个方向映射

还是以鸢尾花数据矩阵 X 为例，矩阵乘法 $X[v_1, v_2]$ 代表 X 将朝着 $[v_1, v_2]$ 两个方向映射。如果 $[v_1, v_2]$ 的取值如图 31 所示，矩阵乘法 $X[v_1, v_2]$ 提取 X 的第 1、3 两列，并将两者顺序调换。

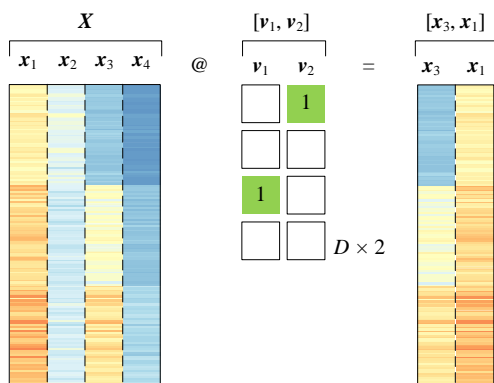


图 31. X 朝两个方向映射

想象一个由鸢尾花四个维度构造的空间 \mathbb{R}^4 ，图 31 相当于将鸢尾花数据映射在一个平面上 \mathbb{R}^2 ，得到的是平面散点图，过程如图 32 所示。

看到这里，大家是否想到了本书第 1 章的成对散点图？每幅散点图的背后实际上都有类似图 31 的矩阵乘法运算。

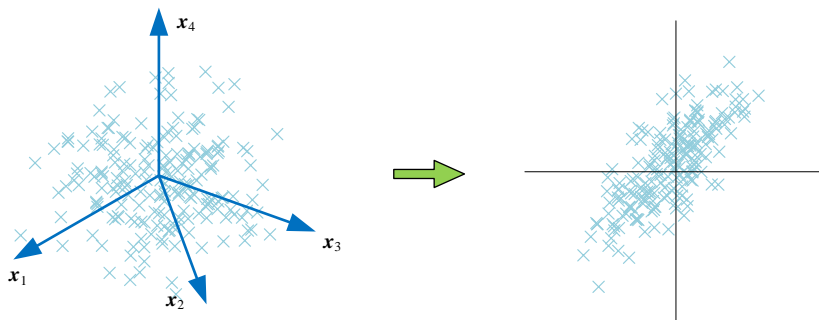


图 32.4 维空间散点压缩到平面上

多个方向映射

矩阵 X 有 D 个维度，可以通过矩阵乘法，将 X 映射到另外一个 D 维度的空间中。

下例中， $V = [v_1, v_2, \dots, v_D]$ ， Z 对应的每一行元素则是新坐标系的坐标值：

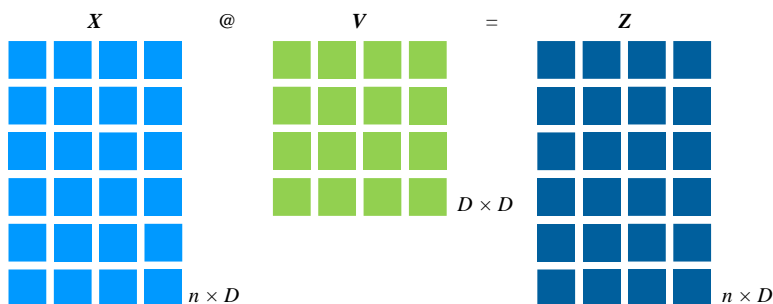
$$XV = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix} \begin{bmatrix} v_1 & v_2 & \dots & v_D \end{bmatrix} = \begin{bmatrix} x^{(1)}v_1 & x^{(1)}v_2 & \dots & x^{(1)}v_D \\ x^{(2)}v_1 & x^{(2)}v_2 & \dots & x^{(2)}v_D \\ \vdots & \vdots & \dots & \vdots \\ x^{(n)}v_1 & x^{(n)}v_2 & \dots & x^{(n)}v_D \end{bmatrix} = Z = \begin{bmatrix} z^{(1)} \\ z^{(2)} \\ \vdots \\ z^{(n)} \end{bmatrix} \quad (81)$$

其中，矩阵 V 为方阵。

如果 V 可逆， V 就是 X 和 Z 相互转化的桥梁：

$$X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix} \xrightleftharpoons[V^{-1}]{V} Z = \begin{bmatrix} z^{(1)} \\ z^{(2)} \\ \vdots \\ z^{(n)} \end{bmatrix} \quad (82)$$

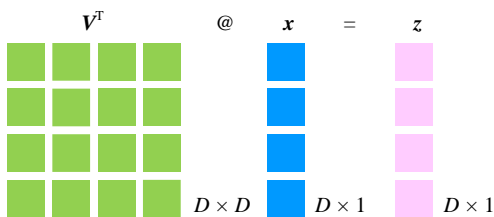
本书第 10 章还会深入讨论上式。

图 33. 一个 D 维度空间 X 数据映射到另一个 D 维度空间

列向量形式

大家见到如下形式时，也不用慌张。如图 34 所示，这个也是上文介绍的映射，只不过 x 为列向量：

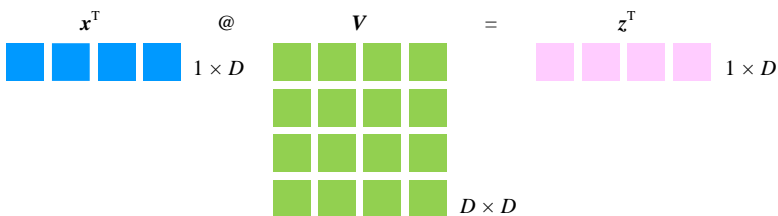
$$V^T x = z \quad (83)$$

图 34. $V^T x = z$ 运算示意图

如图 35 所示，(83) 左右转置，便得到类似 (81) 的结构：

$$x^T V = z^T \quad (84)$$

约定俗成，各种线性代数工具定义偏好列向量；但是，在实际应用中，更常用行向量代表数据点。两者之间的桥梁就是——转置。

图 35. 等式 $V^T x = z$ 左右转置

可以说，本书后续介绍的内容几乎都离不开映射，比如几何变换、正交投影、特征值分解、奇异值分解等等。

5.11 长方阵：奇异值分解、格拉姆矩阵、张量积

本节介绍和长方形矩阵有关的重要矩阵乘法。

奇异值分解

请读者格外注意图 36 所示的矩阵乘法结构。这两种形式经常出现在**奇异值分解** (singular vector decomposition, SVD) 和**主成分分析** (principal component analysis, PCA)。

▲请大家注意图 36 中， D 和 p 的大小关系；不同大小关系对应着不同类型的奇异值分解。本书第 16 章将深入讲解。

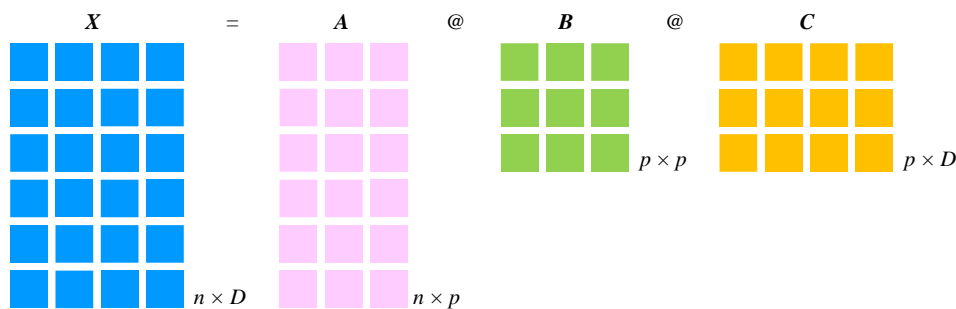


图 36. 三个矩阵相乘

格拉姆矩阵

将矩阵 X 写成一组列向量，如下：

$$X_{n \times D} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,D} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,D} \end{bmatrix} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_D] \quad (85)$$

如图 37 所示，利用 (85)，转置 X^T ($D \times n$) 乘矩阵 X ($n \times D$)，得到一个 $D \times D$ 方阵 $X^T X$ ，可以写成：

$$G = X^T X = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_D^T \end{bmatrix} [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_D] = \begin{bmatrix} \mathbf{x}_1^T \mathbf{x}_1 & \mathbf{x}_1^T \mathbf{x}_2 & \cdots & \mathbf{x}_1^T \mathbf{x}_D \\ \mathbf{x}_2^T \mathbf{x}_1 & \mathbf{x}_2^T \mathbf{x}_2 & \cdots & \mathbf{x}_2^T \mathbf{x}_D \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_D^T \mathbf{x}_1 & \mathbf{x}_D^T \mathbf{x}_2 & \cdots & \mathbf{x}_D^T \mathbf{x}_D \end{bmatrix} \quad (86)$$

上式是矩阵乘法的第一视角。

(86) 中的 G 有自己的名字——**格拉姆矩阵** (Gram matrix)。格拉姆矩阵在数据分析、机器学习算法中有重要作用。

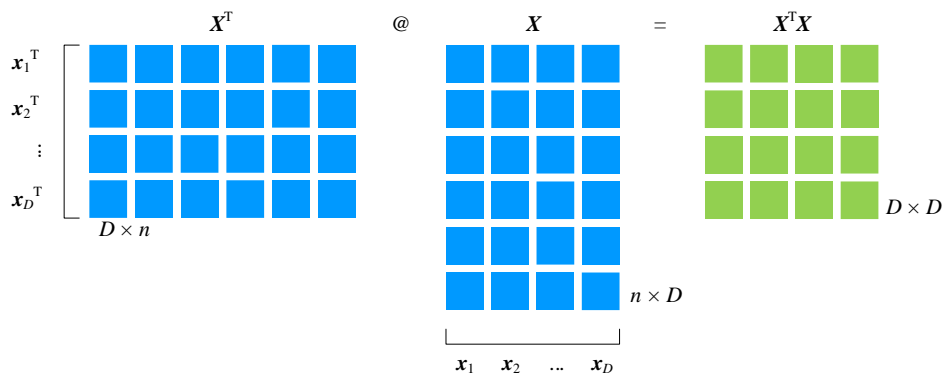


图 37. $X^T X$ 运算过程

如图 38 所示, $X^T X$ 的 (i, j) 元素是 X 中第 i 列向量转置乘以 X 的第 j 列向量:

$$(X^T X)_{i,j} = \mathbf{x}_i^T \mathbf{x}_j \quad (87)$$

当 $i = j$ 时, $\mathbf{x}_i^T \mathbf{x}_i$ 对应的是格拉姆矩阵 G 的对角线元素, 也可以写成 L^2 范数形式 $\|\mathbf{x}_i\|_2^2$ 。

再次强调, 凡是看到矩阵乘积为标量的情况, 要停下来思考一下, 能够将矩阵乘积写成 L^2 范数形式。原因很简单, L^2 范数代表欧氏距离, 这给我们提供一个几何视角。

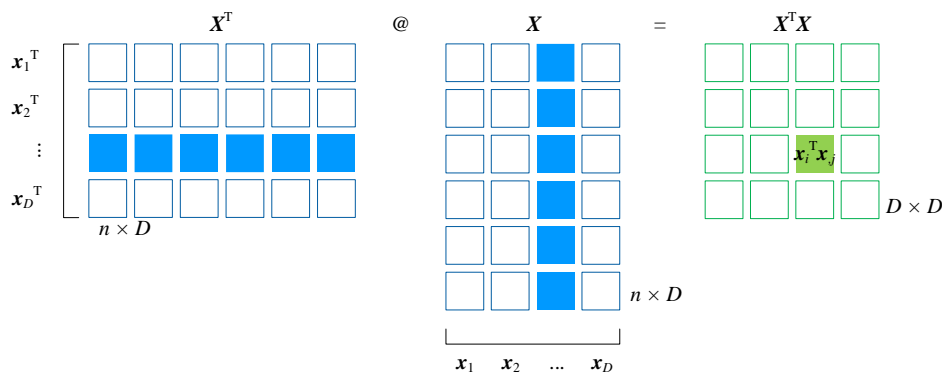


图 38. $X^T X$ 的 (i, j) 元素

标量积

G 还可以写成标量积:

$$\mathbf{G} = \begin{bmatrix} \mathbf{x}_1 \cdot \mathbf{x}_1 & \mathbf{x}_1 \cdot \mathbf{x}_2 & \cdots & \mathbf{x}_1 \cdot \mathbf{x}_D \\ \mathbf{x}_2 \cdot \mathbf{x}_1 & \mathbf{x}_2 \cdot \mathbf{x}_2 & \cdots & \mathbf{x}_2 \cdot \mathbf{x}_D \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_D \cdot \mathbf{x}_1 & \mathbf{x}_D \cdot \mathbf{x}_2 & \cdots & \mathbf{x}_D \cdot \mathbf{x}_D \end{bmatrix} = \begin{bmatrix} \langle \mathbf{x}_1, \mathbf{x}_1 \rangle & \langle \mathbf{x}_1, \mathbf{x}_2 \rangle & \cdots & \langle \mathbf{x}_1, \mathbf{x}_D \rangle \\ \langle \mathbf{x}_2, \mathbf{x}_1 \rangle & \langle \mathbf{x}_2, \mathbf{x}_2 \rangle & \cdots & \langle \mathbf{x}_2, \mathbf{x}_D \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle \mathbf{x}_D, \mathbf{x}_1 \rangle & \langle \mathbf{x}_D, \mathbf{x}_2 \rangle & \cdots & \langle \mathbf{x}_D, \mathbf{x}_D \rangle \end{bmatrix} \quad (88)$$

任何一个单独向量，它的 L^p 范数，特别是 L^2 范数，代表它的“长度”；而几个向量之间的相对关系，则可以通过向量内积来呈现。再进一步，为了方便比较，我们可以用向量夹角余弦值作为度量向量之间相对夹角的数学工具。

格拉姆矩阵之所以重要，一方面是因为它集成了向量长度 (L^2 范数) 和相对夹角 (夹角余弦值) 两部分重要信息。另一方面，格拉姆矩阵 \mathbf{G} 为对称矩阵：

$$\mathbf{G}^T = (\mathbf{X}^T \mathbf{X})^T = \mathbf{X}^T \mathbf{X} = \mathbf{G} \quad (89)$$

一般情况，数据矩阵 \mathbf{X} 都是“细高”长方形矩阵，矩阵运算时这种形状不够友好。比如，细高的 \mathbf{X} 显然不存在转置。而把 \mathbf{X} 转化为方阵 $\mathbf{G} (= \mathbf{X}^T \mathbf{X})$ 之后，很多运算都变得更加容易。

此外， $\mathbf{X}^T \mathbf{X}$ 相当于 \mathbf{X} 的“平方”。大家需要注意 $\mathbf{X}^T \mathbf{X}$ 的单位。比如，鸢尾花数据 \mathbf{X} 单位为厘米， $\mathbf{X}^T \mathbf{X}$ 中每个元素的单位就变成了平方厘米。实践中，碰到矩阵乘法运算，要留意每个矩阵的单位。



本书第 22 章介绍协方差矩阵 (covariance matrix) 时，也将采用类似 (86) 的计算思路。

张量积

将矩阵 \mathbf{X} 写成一系列行向量，如下：

$$\mathbf{X}_{n \times D} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,D} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,D} \end{bmatrix} = \begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \\ \vdots \\ \mathbf{x}^{(n)} \end{bmatrix} \quad (90)$$

利用 (90)，格拉姆矩阵 $\mathbf{X}^T \mathbf{X}$ ，可以写成一系列张量积的和：

$$\mathbf{G} = \mathbf{X}^T \mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)T} & \mathbf{x}^{(2)T} & \cdots & \mathbf{x}^{(n)T} \end{bmatrix} \begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \\ \vdots \\ \mathbf{x}^{(n)} \end{bmatrix} = \sum_{i=1}^n \mathbf{x}^{(i)T} \mathbf{x}^{(i)} = \sum_{i=1}^n \mathbf{x}^{(i)} \otimes \mathbf{x}^{(i)} \quad (91)$$

上式是矩阵乘法的第二视角。

另一个格拉姆矩阵

本节前文的数据矩阵 \mathbf{X} 是细高的，它转置之后得到宽矮的矩阵 \mathbf{X}^T 。而 \mathbf{X}^T 也有自己的格拉姆矩阵，即矩阵 \mathbf{X} ($n \times D$) 乘其转置 \mathbf{X}^T ($D \times n$)，得到一个 $n \times n$ 格拉姆矩阵：

$$\begin{aligned} \mathbf{X}\mathbf{X}^T &= [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_D] \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_D^T \end{bmatrix} \\ &= \mathbf{x}_1\mathbf{x}_1^T + \mathbf{x}_2\mathbf{x}_2^T + \cdots + \mathbf{x}_D\mathbf{x}_D^T \\ &= \sum_{i=1}^D \mathbf{x}_i\mathbf{x}_i^T \end{aligned} \quad (92)$$

观察 (92)，大家是否也发现了张量积的影子？

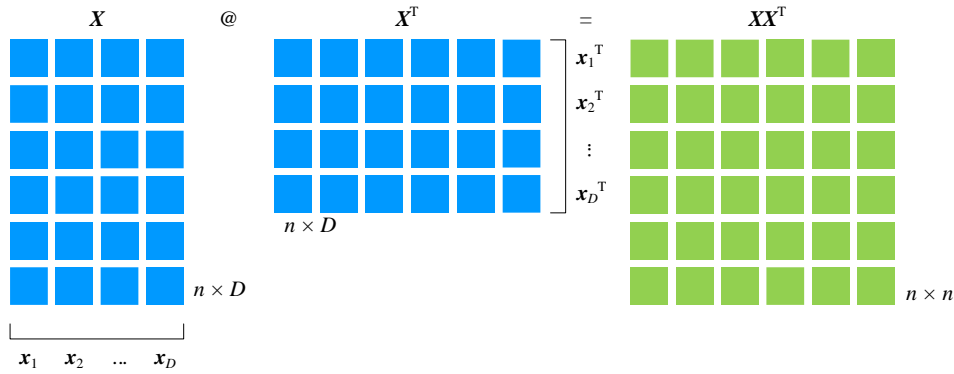


图 39. $\mathbf{X}\mathbf{X}^T$ 运算过程

元素平方和

此外，下式可以计算得到矩阵 \mathbf{X} 的所有元素的平方和：

$$\begin{aligned} \text{trace}(\mathbf{X}^T\mathbf{X}) &= \text{trace} \begin{bmatrix} \mathbf{x}_1 \cdot \mathbf{x}_1 & \mathbf{x}_1 \cdot \mathbf{x}_2 & \cdots & \mathbf{x}_1 \cdot \mathbf{x}_D \\ \mathbf{x}_2 \cdot \mathbf{x}_1 & \mathbf{x}_2 \cdot \mathbf{x}_2 & \cdots & \mathbf{x}_2 \cdot \mathbf{x}_D \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_D \cdot \mathbf{x}_1 & \mathbf{x}_D \cdot \mathbf{x}_2 & \cdots & \mathbf{x}_D \cdot \mathbf{x}_D \end{bmatrix} \\ &= \mathbf{x}_1 \cdot \mathbf{x}_1 + \mathbf{x}_2 \cdot \mathbf{x}_2 + \cdots + \mathbf{x}_D \cdot \mathbf{x}_D \\ &= \sum_{i=1}^n x_{i,1}^2 + \sum_{i=1}^n x_{i,2}^2 + \cdots + \sum_{i=1}^n x_{i,D}^2 = \sum_{j=1}^D \sum_{i=1}^n x_{i,j}^2 \end{aligned} \quad (93)$$

上一章讲解矩阵迹 (trace) 时提到，如果 \mathbf{AB} 和 \mathbf{BA} 都存在， $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$ 。也就是说，对于 (93)，下式成立：

$$\text{trace}(\mathbf{X}^T\mathbf{X}) = \text{trace}(\mathbf{X}\mathbf{X}^T) = \sum_{j=1}^D \sum_{i=1}^n x_{i,j}^2 \quad (94)$$

本书后文还会在不同位置用到 $\text{tr}(AB) = \text{tr}(BA)$ ，请大家格外注意。

此外，向量的范数度量了向量的大小。任意向量 L^2 范数的平方值，就是向量每个元素平方之和。而矩阵 X 的所有元素的平方和实际上也度量了某种矩阵“大小”。一个矩阵的所有元素平方和、再开方叫做矩阵 F -范数。本书第 18 章将介绍常见矩阵范数。

5.12 爱因斯坦求和约定

本书之前的所有矩阵运算都是适用于二阶情况，比如 $n \times D$ 的这种 n 行、 D 列形式。在数据科学和机器学习很多实践中，我们不可避免地要处理高阶矩阵，比如图 40 所示三阶矩阵。Python 中 Xarray 专门用来存储和运算高阶矩阵。

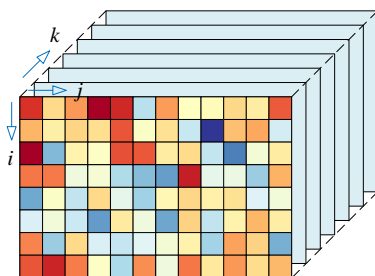


图 40. 三维数组，三阶矩阵

本节则要引出一种简洁表达高阶矩阵运算的数学工具——**爱因斯坦求和约定** (Einstein summation convention 或 Einstein notation)。本节特别介绍如何用 `numpy.einsum()` 函数完成本书前文介绍的主要线性代数运算。此外，PyTorch 中 `torch.einsum()` 函数原理和 `numpy.einsum()` 基本相同，本书不特别介绍。

使用 `numpy.einsum()` 时，大家记住一个要点——输入中重复的索引代表元素相乘，输出中消去的索引意味着相加。

举个例子，矩阵 A 和 B 相乘用 `numpy.einsum()` 函数可以写成：

```
np.einsum('ij,jk->ik', A, B)
```

“->”之前分别为矩阵 A 和 B 的索引，它们用逗号隔开。矩阵 A 行索引为 i ，列索引为 j 。矩阵 B 行索引为 j ，列索引为 k 。 j 为重复索引，因此在这个方向上元素相乘。

“->”之后为输出结果的索引。输出结果索引为 ik ，没有 j ，因此在 j 索引方向上存在求和运算。

表 1 总结如何使用 `numpy.einsum()` 完成常见线性代数运算。现在不需要大家掌握 `numpy.einsum()`。希望大家在日后用到爱因斯坦求和约定时，再回过头来深入学习。

表 1. 使用 `numpy.einsum()` 完成常见线性代数运算

运算	使用 <code>numpy.einsum()</code> 完成运算
向量 a 所有元素求和 (结果为标量)	<code>np.einsum('ij->', a)</code> <code>np.einsum('i->', a_1D)</code>
等行数列向量 a 和 b 的逐项积	<code>np.einsum('ij,ij->ij', a, b)</code> <code>np.einsum('i,i->i', a_1D, b_1D)</code>
等行数列向量 a 和 b 的向量内积 (结果为标量)	<code>np.einsum('ij,ij->', a, b)</code> <code>np.einsum('i,i->', a_1D, b_1D)</code>
向量 a 和自身的张量积	<code>np.einsum('ij,ji->ij', a, a)</code> <code>np.einsum('i,j->ij', a_1D, a_1D)</code>
向量 a 和 b 的张量积	<code>np.einsum('ij,ji->ij', a, b)</code> <code>np.einsum('i,j->ij', a_1D, b_1D)</code>
矩阵 A 的转置	<code>np.einsum('ji', A)</code> <code>np.einsum('ij->ji', A)</code>
矩阵 A 所有元素求和 (结果为标量)	<code>np.einsum('ij->', A)</code>
矩阵 A 对每一列元素求和	<code>np.einsum('ij->j', A)</code>
矩阵 A 对每一行元素求和	<code>np.einsum('ij->i', A)</code>
提取方阵 A 的对角元素 (结果为标量)	<code>np.einsum('ii->i', A)</code>
计算方阵 A 的迹 <code>trace(A)</code> (结果为标量)	<code>np.einsum('ii->', A)</code>
计算矩阵 A 和 B 乘积	<code>np.einsum('ij,jk->ik', A, B)</code>
乘积 AB 结果所有元素求和 (结果为标量)	<code>np.einsum('ij,jk->', A, B)</code>
矩阵 A 和 B 相乘后再转置, 即 $(AB)^T$	<code>np.einsum('ij,jk->ki', A, B)</code>
形状相同矩阵 A 和 B 逐项积	<code>np.einsum('ij,ij->ij', A, B)</code>



表 1 中变量定义和运算都在 `Bk4_Ch5_01.py` 中。

5.13 矩阵乘法的几个雷区

本章最后聊聊运用矩阵乘法时几个潜伏的雷区。

本 PDF 文件为作者草稿，发布目的为方便读者在移动终端学习，终稿内容以清华大学出版社纸质出版物为准。
版权归清华大学出版社所有，请勿商用，引用请注明出处。
代码及 PDF 文件下载：<https://github.com/Visualize-ML>
本书配套微课视频均发布在 B 站——生姜 DrGinger：<https://space.bilibili.com/513194466>
欢迎大家批评指教，本书专属邮箱：jiang.visualize.ml@gmail.com

不满足交换律

代数中，乘法满足交换律，比如 $ab = ba$ 。

但是，一般情况，矩阵乘法不满足交换律：

$$\mathbf{AB} \neq \mathbf{BA} \quad (95)$$



本书在第 8 章将通过几何变换角度解释为什么矩阵乘法一般不满足交换律。

平方

如果方阵 \mathbf{A} 和 \mathbf{B} 满足：

$$\mathbf{A}^2 = \mathbf{B}^2 \quad (96)$$

不能得到：

$$\mathbf{A} = \pm \mathbf{B} \quad (97)$$

对于非方阵 \mathbf{A} ， $\mathbf{A}^T \mathbf{A}$ 或 $\mathbf{A} \mathbf{A}^T$ 相当于 \mathbf{A} 的“平方”。

如果 \mathbf{A} 和 \mathbf{B} 非方阵，且满足：

$$\mathbf{A}^T \mathbf{A} = \mathbf{B}^T \mathbf{B} \quad (98)$$

上式也无法推导得到 (97)。

同理，下式也无法推导得到 (97)：

$$\mathbf{A} \mathbf{A}^T = \mathbf{B} \mathbf{B}^T \quad (99)$$

和的平方

代数中， $(a + b)^2 = a^2 + 2ab + b^2$ 。

如果 \mathbf{A} 和 \mathbf{B} 为方阵，两者和的平方展开得到：

$$(\mathbf{A} + \mathbf{B})^2 = (\mathbf{A} + \mathbf{B})(\mathbf{A} + \mathbf{B}) = \mathbf{A}^2 + \mathbf{AB} + \mathbf{BA} + \mathbf{B}^2 \quad (100)$$

上式中， \mathbf{AB} 和 \mathbf{BA} 不能随意合并。

如果 \mathbf{A} 和 \mathbf{B} 非方阵，下式相当于两者和的平方：

$$(\mathbf{A} + \mathbf{B})^T (\mathbf{A} + \mathbf{B}) = (\mathbf{A}^T + \mathbf{B}^T)(\mathbf{A} + \mathbf{B}) = \mathbf{A}^T \mathbf{A} + \mathbf{A}^T \mathbf{B} + \mathbf{B}^T \mathbf{A} + \mathbf{B}^T \mathbf{B} \quad (101)$$

上式显然不同于下式：

$$(\mathbf{A} + \mathbf{B})(\mathbf{A} + \mathbf{B})^T = (\mathbf{A} + \mathbf{B})(\mathbf{A}^T + \mathbf{B}^T) = \mathbf{A} \mathbf{A}^T + \mathbf{A} \mathbf{B}^T + \mathbf{B} \mathbf{A}^T + \mathbf{B} \mathbf{B}^T \quad (102)$$

矩阵相等

如果 $a \neq 0$, $ab = ac$ 可以推导 $a(b - c) = 0$, 继而得到 $b = c$ 。但是矩阵乘法中, 如果 \mathbf{A} 不是零矩阵, 即 $\mathbf{A} \neq \mathbf{O}$, 并且:

$$\mathbf{AB} = \mathbf{AC} \quad (103)$$

可以推导得到:

$$\mathbf{A}(\mathbf{B} - \mathbf{C}) = \mathbf{O} \quad (104)$$

但是, 不能直接得出 $\mathbf{B} = \mathbf{C}$ 。

举个例子, 给定:

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 2 & 2 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{C} = \begin{bmatrix} 4 & 2 \\ -1 & 1 \end{bmatrix} \quad (105)$$

如下等式成立:

$$\underbrace{\begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}}_{\mathbf{A}} @ \underbrace{\begin{bmatrix} 2 & 2 \\ 0 & 1 \end{bmatrix}}_{\mathbf{B}} = \underbrace{\begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}}_{\mathbf{A}} @ \underbrace{\begin{bmatrix} 4 & 2 \\ -1 & 1 \end{bmatrix}}_{\mathbf{C}} = \begin{bmatrix} 2 & 4 \\ 4 & 8 \end{bmatrix} \quad (106)$$

显然 $\mathbf{B} \neq \mathbf{C}$ 。这是因为 (105) 给出的矩阵 \mathbf{A} 不可逆!

如果 \mathbf{A} 可逆, 我们需要“老老实实”地在等式 (103) 左右分别左乘 \mathbf{A}^{-1} , 一步步推导得到:

$$\mathbf{A}^{-1}(\mathbf{AB}) = \mathbf{A}^{-1}(\mathbf{AC}) \Rightarrow (\mathbf{A}^{-1}\mathbf{A})\mathbf{B} = (\mathbf{A}^{-1}\mathbf{A})\mathbf{C} \Rightarrow \mathbf{IB} = \mathbf{IC} \Rightarrow \mathbf{B} = \mathbf{C} \quad (107)$$

如果下式对于 \mathbb{R}^n 中任意 \mathbf{x} 都成立, 则 $\mathbf{A} = \mathbf{B}$:

$$\mathbf{A}_{m \times n} \mathbf{x}_{n \times 1} = \mathbf{B}_{m \times n} \mathbf{x}_{n \times 1} \quad (108)$$

零矩阵

如果 $ab = 0$, 可以得知 $a = 0$ 或 $b = 0$ 。但是, 如果 $\mathbf{AB} = \mathbf{O}$, 则无法得到 $\mathbf{A} = \mathbf{O}$ 或 $\mathbf{B} = \mathbf{O}$ 。

举个例子, 如下 \mathbf{A} 和 \mathbf{B} 乘积为零矩阵, 但是显然它们都不是零矩阵。

$$\underbrace{\begin{bmatrix} 0 & 1 \\ 0 & 2 \end{bmatrix}}_{\mathbf{A}} @ \underbrace{\begin{bmatrix} 3 & 4 \\ 0 & 0 \end{bmatrix}}_{\mathbf{B}} = \underbrace{\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}}_{\mathbf{O}} \quad (109)$$

如果 $k\mathbf{A} = \mathbf{O}$, 标量 $k = 0$ 或矩阵 $\mathbf{A} = \mathbf{O}$ 。

注意, $\mathbf{AO} = \mathbf{O}$ 中, 两个零矩阵的形状很可能不一致。

注意顺序

在多个矩阵连乘展开遇到求逆或置换时，大家需要格外注意调换顺序，比如：

$$\begin{aligned}(ABC)^{-1} &= C^{-1}B^{-1}A^{-1} \\ (ABC)^T &= C^TB^TA^T \\ ((ABC)^{-1})^T &= (C^{-1}B^{-1}A^{-1})^T = (A^{-1})^T(B^{-1})^T(C^{-1})^T = (A^T)^{-1}(B^T)^{-1}(C^T)^{-1}\end{aligned}\tag{110}$$

其中， A 、 B 、 C 均可逆。注意， $(A^T)^{-1} = (A^{-1})^T$ 。

标量乘法

遇到标量乘法时注意：

$$\begin{aligned}(kA)^{-1} &= \frac{1}{k}A^{-1} \\ (kA)^T &= kA^T\end{aligned}\tag{111}$$

上式中， k 非零。此外， A^{-1} 代表矩阵的逆，不能类比成代数中的“倒数”。因此， A^{-1} 不能写成 $1/A$ 或 $\frac{1}{A}$ ，它俩在线性代数中没有定义。

结果可能是标量

矩阵的乘积结果可能是个标量。本章前文给出过几个例子，总结如下：

$$\begin{aligned}\mathbf{x} \cdot \mathbf{y} &= \mathbf{y} \cdot \mathbf{x} = \langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = (\mathbf{x}^T \mathbf{y})^T = \mathbf{y}^T \mathbf{x} \\ \mathbf{I} \cdot \mathbf{x} &= \mathbf{I}^T \mathbf{x} = \mathbf{x}^T \mathbf{I} \\ \mathbf{x} \cdot \mathbf{x} &= \langle \mathbf{x}, \mathbf{x} \rangle = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|_2^2 \\ \sqrt{\mathbf{x} \cdot \mathbf{x}} &= \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\mathbf{x}^T \mathbf{x}} = \|\mathbf{x}\|_2 \\ \mathbf{I}^T \mathbf{X} \mathbf{I}, \quad \mathbf{x}^T \mathbf{Q} \mathbf{x}, \quad \mathbf{v}_i^T \Sigma \mathbf{v}_j\end{aligned}\tag{112}$$

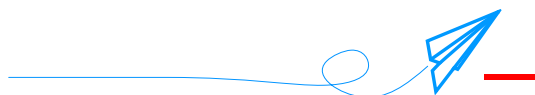
上述结果相当于是 1×1 矩阵。 1×1 矩阵的转置为其本身。反复强调，遇到矩阵乘积为标量的情况，请大家考虑矩阵乘积能否看做是某种“距离”。

不能消去

当矩阵乘积为标量时，写在算式分母上就不足为奇了，比如：

$$\frac{\mathbf{x}^T \mathbf{y}}{\mathbf{x}^T \mathbf{x}}, \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}, \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{B} \mathbf{x}}, \frac{\mathbf{x}^T \mathbf{A}^T \mathbf{Q} \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}} \quad (113)$$

如果分子、分母上都出现同一个矩阵，绝不能消去。显然，(113) 中 \mathbf{x}^T 和 \mathbf{x} 都不能消去！最后一个分式中 \mathbf{A}^T 和 \mathbf{A} 也不能消去。



本章全景展示了常见的矩阵乘法形态。每种形态的矩阵乘法都很重要，因此也不能用四幅图来总结本章主要内容。

大家想要活用线性代数这个宝库中的各种数学工具，熟练掌握矩阵乘法规则是绕不过去的一道门槛。再强调一次，大家在学习不同的矩阵运算时，要试图从几何和数据这两个视角去理解。这也是本书要特别强化的一点。

此外，本章针对矩阵乘法运算也没有给出任何代码，因为在 Numpy 中矩阵乘法常用运算符就是 @。本章还介绍了爱因斯坦求和约定，不要求大家掌握。本章最后聊了聊矩阵乘法运算的常见雷区，希望大家格外小心。

矩阵乘法规则像是枷锁，它条条框框、冷酷无情、不容妥协；但是，在枷锁下，我们看到了矩阵乘法的另一面——无拘无束、血脉偾张、海纳百川。

希望大家一边学习本书剩余内容，一边能够不断回头看这一章内容，相信大家一定会和我一样，叹服于矩阵乘法展现出来的自由、包容，和纯粹的美。