

23

Four Vector Spaces

数据空间

用 SVD 分解寻找数据矩阵的四个空间



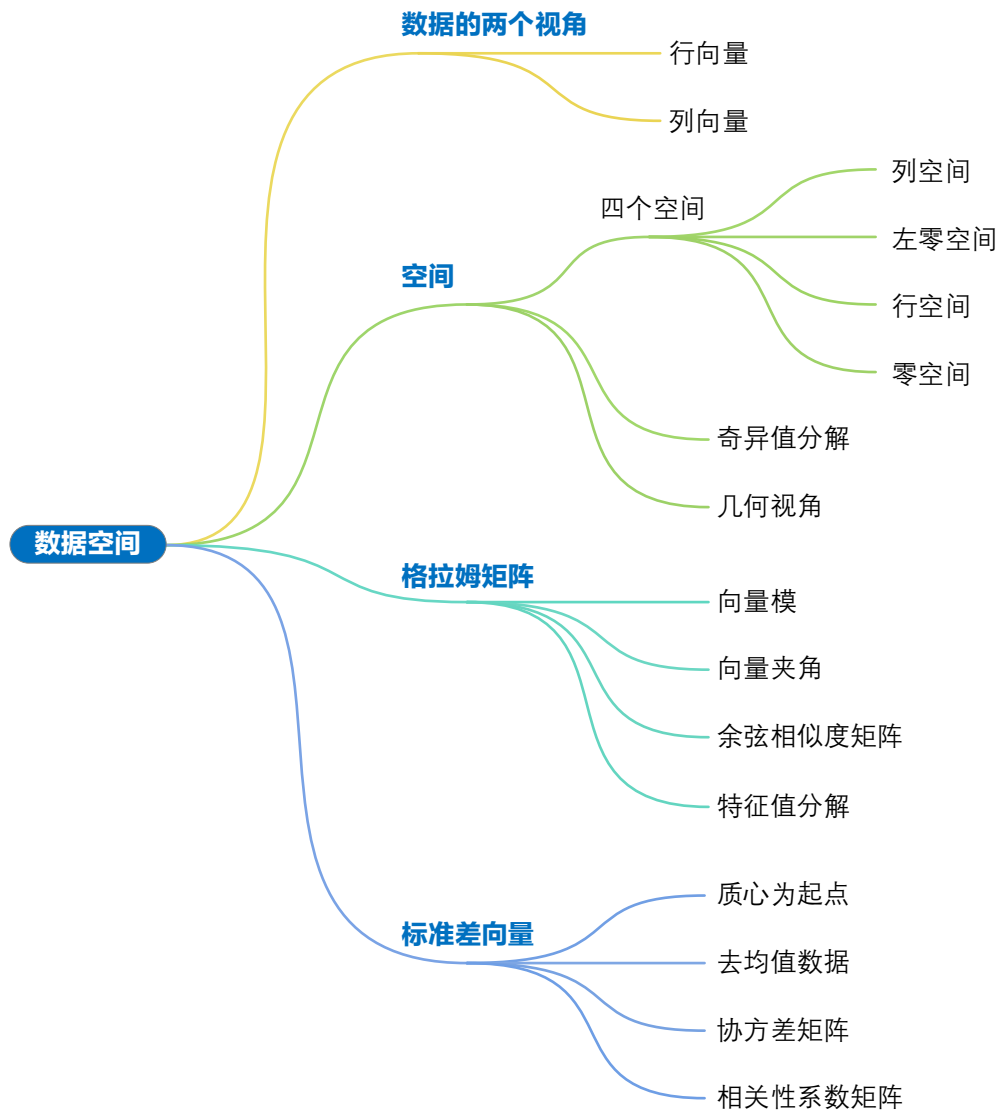
智慧的真正标志不是知识，而是想象力。

The true sign of intelligence is not knowledge but imagination.

—— 阿尔伯特·爱因斯坦 (Albert Einstein) | 理论物理学家 | 1879 ~ 1955



- numpy.cov() 计算协方差矩阵
- numpy.corr() 计算相关性系数矩阵
- numpy.diag() 如果 A 为方阵，numpy.diag(A) 函数提取对角线元素，以向量形式输入结果；如果 a 为向量，numpy.diag(a) 函数将向量展开成方阵，方阵对角线元素为 a 向量元素
- numpy.linalg.eig() 特征值分解
- numpy.linalg.inv() 计算逆矩阵
- numpy.linalg.norm() 计算范数
- seaborn.heatmap() 绘制热图



23.1 从数据矩阵 X 说起

本书最后三章叫“数据三部曲”，这三章一方面从数据、空间、几何角度总结全书前文核心内容，另外一方面介绍这些数学工具在数据科学和机器学习领域的应用。

毫不夸张地说，没有线性代数就没有现代计算，大家将会在本系列丛书《数据科学》和《机器学习》两册书的每个角落看到矩阵运算。

“多重视角”仍然是这三章的特色。线性代数中向量、空间、投影、矩阵、矩阵分解等数学工具天然地弥合代数、几何、数据之间的鸿沟。

本章是“数据三部曲”的第一章，将以数据矩阵为切入点，主要通过奇异值分解和大家探讨四个重要的空间定义和用途。

数据矩阵

数据矩阵 (data matrix) 不过就是以表格形式存储的数据。

除了表格功能，矩阵更重要的功能是——**线性映射** (linear mapping)。而矩阵乘法是线性映射的核心。矩阵分解不过是矩阵连乘，将一个复杂的几何变换拆解成容易理解的成分，比如缩放、旋转、投影、剪切等等。

本书最开始便介绍过，数据矩阵可以从两个角度观察。数据矩阵 X 的每一行是一个行向量，代表一个样本观察值； X 的每一列为一个列向量，代表某个特征上的样本数据。

行向量

回顾前文，为了区分数据矩阵中的行向量和列向量，本书中数据矩阵的行向量序号采用上标加括号记法，比如：

$$X_{n \times D} = \begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \\ \vdots \\ \mathbf{x}^{(n)} \end{bmatrix} \quad (1)$$

其中，第 i 行行向量 D 个元素为：

$$\mathbf{x}^{(i)} = [x_{i,1} \quad x_{i,2} \quad \cdots \quad x_{i,D}] \quad (2)$$

图 1 所示为从行向量角度观察数据矩阵，每一个行向量 $\mathbf{x}^{(i)}$ 代表坐标系中一个点。所有数据散点构成坐标系中的“云”。

实际上，行向量也是具有方向和大小的向量，也可以看成是箭头，因此也有自己的空间。这是本书马上要探讨的内容。

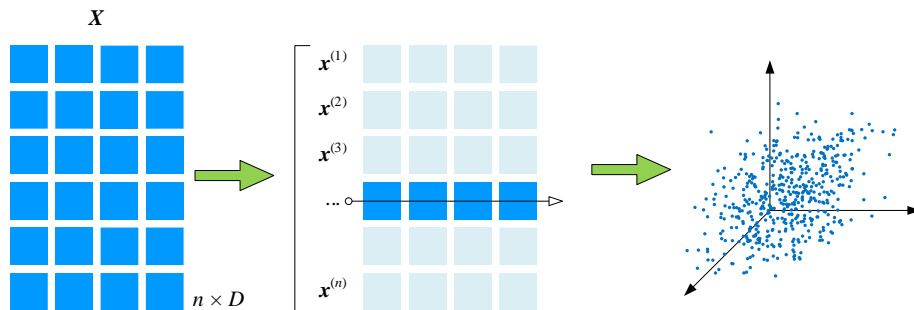


图 1. 从行向量角度观察数据矩阵

列向量

数据矩阵的列向量序号采用下标记法，比如：

$$\mathbf{X}_{n \times D} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_D] \quad (3)$$

其中，第 j 列列向量 n 个元素为：

$$\mathbf{x}_j = \begin{bmatrix} x_{1,j} \\ x_{2,j} \\ \vdots \\ x_{n,j} \end{bmatrix} \quad (4)$$

如图 2 所示，从几何角度，数据矩阵 \mathbf{X} 的所有列向量 (蓝色箭头) 的起始点均在原点 $\mathbf{0}$ 。 $[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_D]$ 这些向量的长度和方向信息均包含在格拉姆矩阵 $\mathbf{G} = \mathbf{X}^T \mathbf{X}$ 之中。

向量长度的表现形式为向量的模，即 L^2 范数。

向量方向是两两向量之间的相对夹角。更具体地说，是两两向量夹角余弦值。

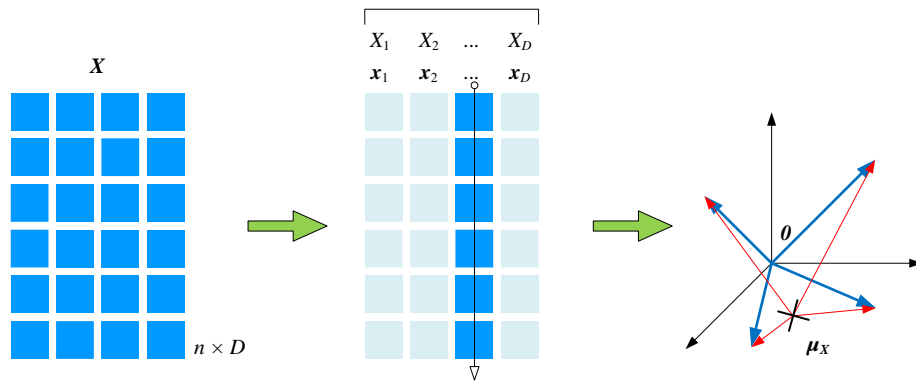


图 2. 从列向量角度观察数据矩阵

如果将图 2 向量起点移动到数据质心 μ_X (即 $E(X)$ 的转置), 这时向量 (红色箭头) 的长度可以看做是标准差 (的若干倍), 而向量之间夹角为随机变量之间的线性相关系数。

从统计角度来看, 将向量起点移动到 μ_X 实际上就是数据矩阵 X 去均值, 即中心化, 对应运算为 $X_c = X - E(X)$ 。本章后文还将深入介绍这一重要视角。

协方差矩阵 Σ 相当于是 X_c 的格拉姆矩阵。准确来说, 对于样本数据, $X_c^T X_c = (n-1)\Sigma$ 。协方差矩阵 Σ 包含了样本标准差和线性相关系数等信息。

区分符号

现在有必要再次强调本系列丛书的容易混淆的代数、线性代数和概率统计符号。

粗体、斜体、小写 x 为列向量。从概率统计的角度, x 可以代表随机变量 X 采样得到的样本数据, 偶尔也代表 X 总体数据。随机变量 X 样本数据集合为 $X = \{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$ 。

粗体、斜体、小写、加下标序号的 x_1 为列向量, 下角标仅仅是序号, 以便区分 x_1 、 x_2 、 x_j 、 x_D 等等。从概率统计的角度, x_1 可以代表随机变量 X_1 样本数据, 也可以表达 X_1 总体数据。

行向量 $x^{(1)}$ 代表一个具有多个特征的样本点。

从代数角度, 斜体、小写、非粗体 x_1 代表变量, 下角标代表变量序号。这种记法常用在函数解析式中, 比如线性回归解析式 $y = x_1 + x_2$ 。

$x^{(1)}$ 代表变量 x 的一个取值, 或代表随机变量 X 的一个取值。

而 $x_1^{(1)}$ 代表变量 x_1 的一个取值, 或代表随机变量 X_1 的一个取值, 比如 $X_1 = \{x_1^{(1)}, x_1^{(2)}, \dots, x_1^{(n)}\}$ 。

粗体、斜体、大写 X 则专门用来表达多行、多列的数据矩阵, $X = [x_1, x_2, \dots, x_D]$ 。数据矩阵 X 中第 i 行、第 j 列元素则记做 x_{ij} 。多元线性回归中, X 也叫**设计矩阵** (design matrix)。

我们还会用粗体、斜体、小写希腊字母 χ (chi, 读作/'kaɪ/) 代表 D 维随机变量构成的列向量, $\chi = [X_1, X_2, \dots, X_D]^T$ 。希腊字母 χ 主要用在多元概率统计中。

23.2 向量空间：从 SVD 分解角度理解

这一节介绍 X 列向量和行向量张成的四个空间以及它们之间关系。

列向量：列空间、左零空间

由 X 的列向量 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_j, \dots, \mathbf{x}_D$ 张成的子空间 $\text{span}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_D)$ 为 X 的列空间 (column space)，记做 $C(X)$ 。很多书上也把列空间记做 $\text{Col}(X)$ ，或 $\text{Col}X$ 。

与 $C(X)$ 相对应的是左零空间 (left null space)，记做 $\text{Null}(X^T)$ 。 $C(X)$ 和 $\text{Null}(X^T)$ 构成了 \mathbb{R}^n 。 X 的列向量元素个数为 n ，因此需要匹配空间 \mathbb{R}^n ，才能“装下” X 的列向量。

而 $C(X)$ 和 $\text{Null}(X^T)$ 分别都是 \mathbb{R}^n 的子空间，两者的维度之和为 n ，即 $\dim(C(X)) + \dim(\text{Null}(X^T)) = n$ 。

$C(X)$ 和 $\text{Null}(X^T)$ 互为正交补 (orthogonal complement)，即：

$$C(X)^\perp = \text{Null}(X^T) \quad (5)$$

行向量：行空间、零空间

由 X 的行向量 $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(j)}, \dots, \mathbf{x}^{(n)}$ 张成的子空间 $\text{span}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)})$ 为 X 的行空间 (row space)，记做 $R(X)$ 。很多书上也记做 $\text{Row}(X)$ 或 $\text{Row}X$ 。

与 $R(X)$ 相对应的是零空间 (null space)，也叫右零空间 (right null space)，记做 $\text{Null}(X)$ 。

X 的行向量元素数量为 D ，空间 \mathbb{R}^D 才能“装下” X 的行向量。 $R(X)$ 和 $\text{Null}(X)$ 构成了 \mathbb{R}^D 。 $R(X)$ 和 $\text{Null}(X)$ 分别都是 \mathbb{R}^D 的子空间。

$R(X)$ 和 $\text{Null}(X)$ 互为正交补，即：

$$R(X)^\perp = \text{Null}(X) \quad (6)$$

$R(X)$ 的维度为 $\dim(R(X)) = \text{rank}(X)$ 。 $R(X)$ 和 $\text{Null}(X)$ 的维度之和为 D ，即 $\dim(R(X)) + \dim(\text{Null}(X)) = D$ 。也就是说，只有 X 非满秩， $\text{Null}(X)$ 维数才不为 0。

怎么理解这四个空间？

相信大家读完本节前文这四个空间定义已经晕头转向，云里雾里不知所云。

的确，这四个空间的定义让很多人望而却步。很多线性代数教材多是从线性方程组 $Ax = b$ 角度讲解这四个空间，而作者认为这个视角并没有降低理解这四个空间的难度。

下面，我们从数据和几何两个角度来理解这四个空间，并且介绍如何将它们和本书前文介绍的向量内积、格拉姆矩阵、向量空间、子空间、秩、特征值分解、SVD 分解、数据质心、协方差矩阵等线性代数概念联系起来。

从完全型 SVD 分解说起

对“细长”矩阵 X 进行完全型 SVD 分解，得到等式：

$$X = USV^T \quad (7)$$

图 3 所示为 X 完全型 SVD 分解示意图。

⚠ 请大家注意几个矩阵形状。完全型 SVD 分解， X 和 S 为一般为细高型， U 和 V 为方阵。

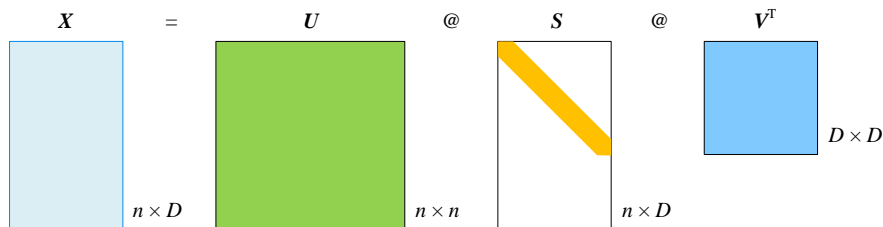


图 3. X 进行完全型 SVD 分解

根据前文所学，大家应该很清楚 U 为 $n \times n$ 正交矩阵，也就是说 U 列向量 $[u_1, u_2, \dots, u_n]$ 特点是两两正交 (向量内积为 0)，且向量模均为 1。

$[u_1, u_2, \dots, u_n]$ 为张成 \mathbb{R}^n 空间的一组规范正交基。

同理， V 为 $D \times D$ 正交矩阵，因此 $V = [v_1, v_2, \dots, v_D]$ 是张成 \mathbb{R}^D 空间的一组规范正交基。

如图 4 所示， U 和 V 之间的联系为 $US = XV$ 。

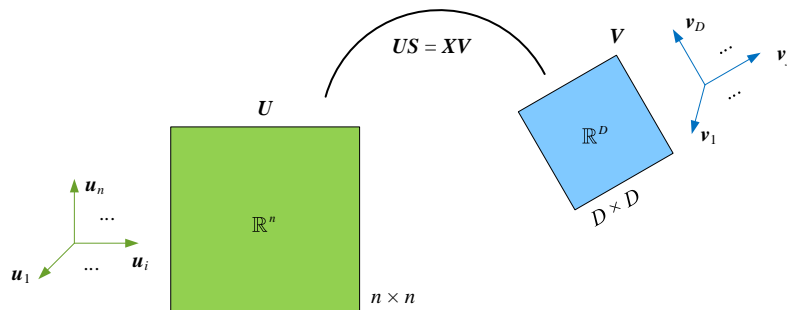


图 4. 对于矩阵 X 来说， \mathbb{R}^n 空间和 \mathbb{R}^D 空间关系

另外，对“粗短” X^T 矩阵进行完全型 SVD 分解，就是对 (7) 转置：

$$X^T = (USV^T)^T = VS^T U^T \quad (8)$$

图 5 所示为 X^T 进行完全型 SVD 分解示意图。后面，我们会用到这一分解。

▲ 注意，对于完全型 SVD 分解，奇异值矩阵 S 虽然是对角阵，但不是方阵，因此 $S^T \neq S$ 。

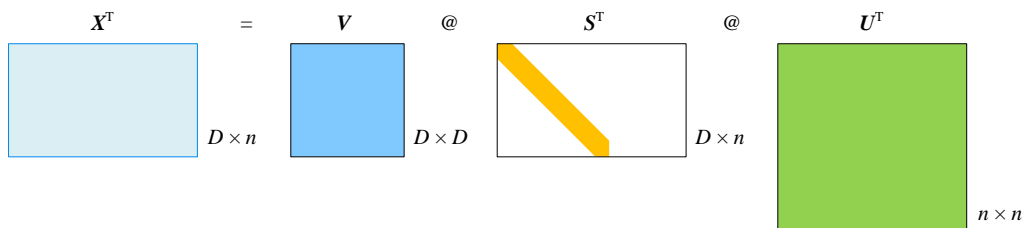


图 5. X^T 进行完全型 SVD 分解

23.3 紧凑型 SVD 分解：剔除零空间

紧凑型 SVD 分解

在讲解奇异值分解时，我们特别介绍了紧凑型 SVD 分解。紧凑型 SVD 分解对应的情况为 $\text{rank}(X) = r < D$ 。奇异值矩阵 S 可以分成四个子块：

$$S = \begin{bmatrix} S_{r \times r} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (9)$$

上式中，矩阵 $S_{r \times r}$ 对角线元素为非 0 奇异值。

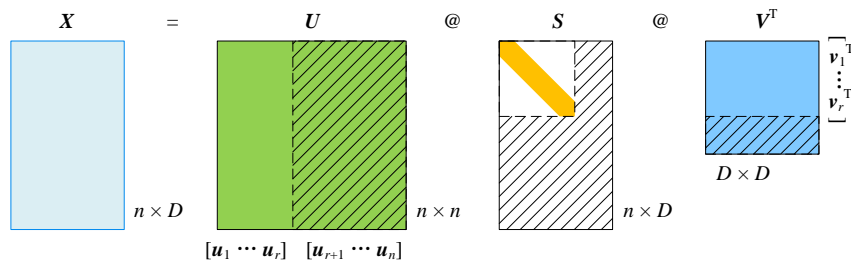
图 6 所示为 X 进行紧凑型 SVD 分解示意图。本书第 16 章介绍过，分块矩阵乘法中，图 6 中阴影部分对应的分块矩阵可以全部消去。

正交矩阵 U 保留 $[u_1, \dots, u_r]$ 子块，消去 $[u_{r+1}, \dots, u_n]$ 。

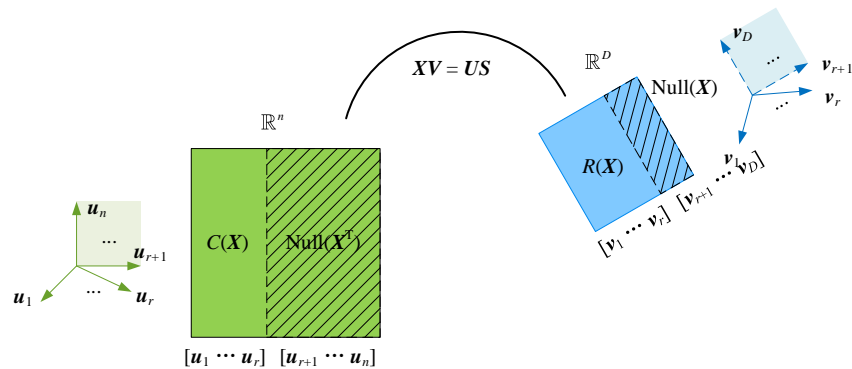
正交矩阵 V 保留 $[v_1, \dots, v_r]$ 子块，消去 $[v_{r+1}, \dots, v_D]$ 。

$[u_1, \dots, u_r]$ 是 X 的列空间 $C(X)$ 基底。而 $[v_1, \dots, v_r]$ 是 X 的行空间 $R(X)$ 基底。

▲ 注意，图 6 中 V 存在转置运算。

图 6. X 进行紧凑型 SVD 分解

实际上, U 和 V 矩阵中消去的子块和上一节说到的**零空间**有直接联系。先给出图 7 这幅图, 我们马上展开讲解。

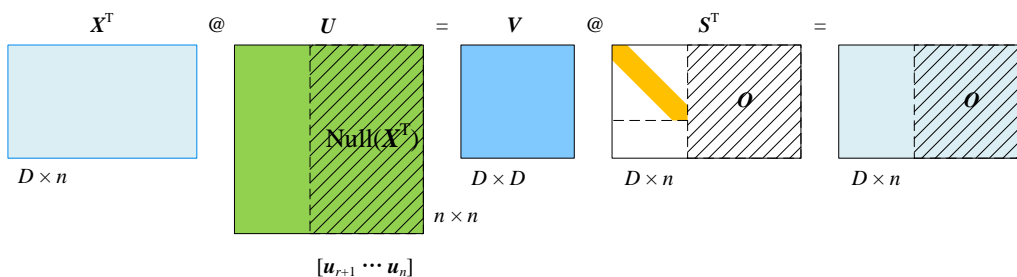
图 7. \mathbb{R}^n 空间和 \mathbb{R}^D 空间关系, 考虑**列空间**、**左零空间**、**行空间**、**零空间**

列空间, 左零空间

$[u_1, \dots, u_r]$ 是 X 的**列空间** $C(X)$ 基底。而 $[u_{r+1}, \dots, u_n]$ 是**左零空间** $\text{Null}(X^T)$ 基底。

如图 8 所示, 将 S^T 左右分块, 右侧分块矩阵为 O 矩阵。 X^T 向**左零空间** $\text{Null}(X^T)$ $[u_{r+1}, \dots, u_n]$ 投影的结果为全 0 矩阵 O 。

白话说, \mathbb{R}^n 用来装 X 的**列向量**, 绝对“杀鸡用牛刀”。 $[u_1, \dots, u_r]$ 张起的子空间就“刚刚好”够装下 X 的**列向量**。而 \mathbb{R}^n 中没有被用到的部分就是 $[u_{r+1}, \dots, u_n]$ 张起的**左零空间** $\text{Null}(X^T)$ 。

图 8. X^T 向 $\text{Null}(X^T)$ $[u_{r+1}, u_2, \dots, u_n]$ 投影的结果为 O

这就是为什么 $\text{Null}(X^T)$ 被称作左“零”空间的原因，因为投影结果为零矩阵。而且，我们也同时在图 8 中投影运算中 X^T 看到了“转置”，这就解释了为什么列空间 $C(X)$ 对应 $\text{Null}(X^T)$ 。

多说一句，(8) 可以写成：

$$X^T U = V S^T \quad (10)$$

上式正交投影中，矩阵 X^T 对应的投影矩阵是 U ， X^T 的每一行代表一个散点，对应 X 的列向量。大家在这句话中看到列空间 $C(X)$ 和 $\text{Null}(X^T)$ 中“列”和“ X^T ”这两个字眼了吧！

行空间，零空间

而 $[v_1, \dots, v_r]$ 是 X 的行空间 $R(X)$ 基底。 $[v_{r+1}, \dots, v_D]$ 是零空间 $\text{Null}(X)$ 的基底。

白话说， \mathbb{R}^D 来装 X 的行向量，可能大材小用，也可能大小合适。 $[v_1, \dots, v_r]$ 张起的子空间就刚刚好够装下 X 的行向量。富余的部分就是 $[v_{r+1}, \dots, v_D]$ 张起的零空间 $\text{Null}(X)$ 。 $\text{rank}(X) = r = D$ 时， \mathbb{R}^D 装 X 的行向量后没有任何余量。

也用正交投影视角来看，将 (8) 写成：

$$X V = U S \quad (11)$$

矩阵 X 对应的投影矩阵是 V ， X 的每一行代表一个散点，对应 X 的行向量。如图 9 所示，将 S 左右分块，右侧分块矩阵为 O 矩阵。 X 向 $\text{Null}(X)$ 投影的结果为 Z 的右侧零矩阵 O 。

图 9 解释了为什么 $\text{Null}(X)$ 被称作“零”空间，而行空间 $R(X)$ 对应零空间 $\text{Null}(X)$ 。

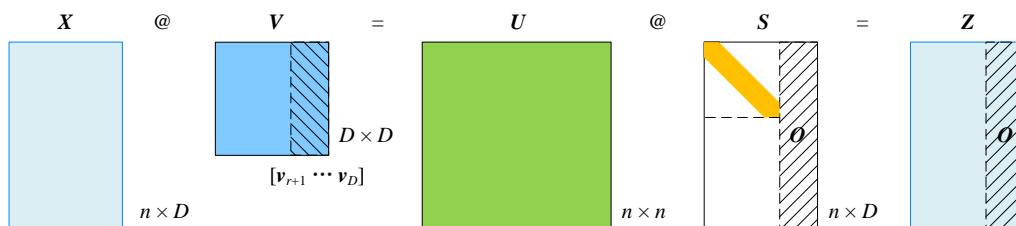


图 9. X 向 $\text{Null}(X)$ 投影的结果为 O

复盘一下

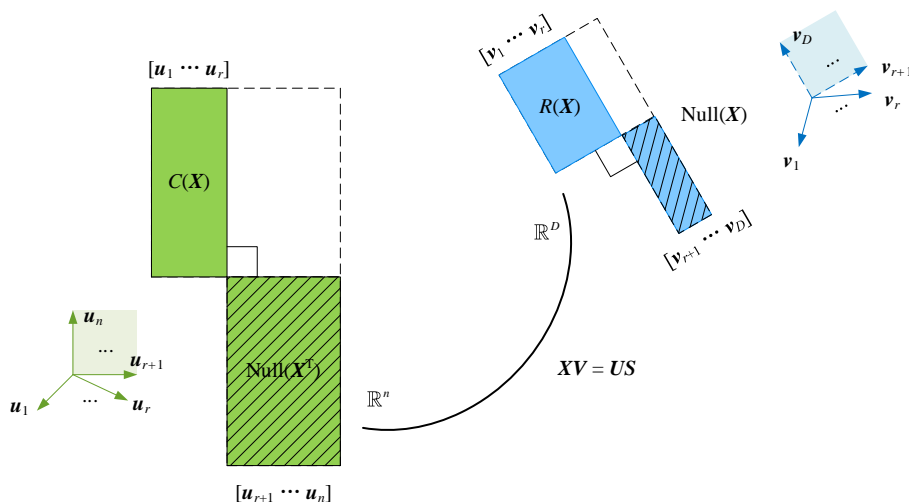
U 是正交矩阵，即 $[u_1, u_2, \dots, u_n]$ 中列向量两两垂直。基底 $[u_1, u_2, \dots, u_n]$ 张起了 \mathbb{R}^n 。

将 $[u_1, u_2, \dots, u_n]$ 划分成两块—— $C(X) = [u_1, \dots, u_r]$ 、 $\text{Null}(X^T) = [u_{r+1}, \dots, u_n]$ 。**列空间** $C(X)$ 和 **左零空间** $\text{Null}(X^T)$ 互为正交补。

“正交”两字，来自于 $[u_1, u_2, \dots, u_n]$ 中列向量两两垂直。“补”字，可以理解为“补齐”，也就是说 $C(X)$ 和 $\text{Null}(X^T)$ 补齐了 \mathbb{R}^n 。

同理，将 $V = [v_1, v_2, \dots, v_D]$ 划分成两块—— $R(X) = [v_1, \dots, v_r]$ 、 $\text{Null}(X) = [v_{r+1}, \dots, v_D]$ 。而**行空间** $R(X)$ 和 **零空间** $\text{Null}(X)$ 互为正交补，两者“补齐”得到 \mathbb{R}^D 。

在图 7 基础上，考虑这两对正交关系，加上 \mathbb{R}^n 空间和 \mathbb{R}^D 空间，我们用图 10 可视化这六个空间。图中加阴影的部分对应**左零空间**和**零空间**。

图 10. \mathbb{R}^n 空间和 \mathbb{R}^D 空间关系，考虑**列空间**、**左零空间**、**行空间**、**零空间**的正交关系

四个空间：因 X 而生

格外强调， \mathbb{R}^n 空间和 \mathbb{R}^D 空间是“永恒”存在的，是“铁打的庙”。但是，能张成这两个空间的规范正交基有无数组，都是“流水的和尚”。

$[u_1, \dots, u_n]$ ，即 $C(X) + \text{Null}(X^T)$ ，是张成 \mathbb{R}^n 空间无数组规范正交基中的一组。

$[v_1, \dots, v_D]$ ，即 $R(X) + \text{Null}(X)$ ，是张成 \mathbb{R}^D 空间无数组规范正交基中的一组。

值得强调的是，在矩阵 X 眼中， $C(X)$ 、 $\text{Null}(X^T)$ 、 $R(X)$ 和 $\text{Null}(X)$ 是独一无二的存在，因为它们都是为矩阵 X 而生！

也就是说，数据矩阵 X 稍有变化，不管是元素、还是形状变化，这四个空间就会随之变化。

而获得这四个空间最便捷的方法就是堪称宇宙第一矩阵分解的奇异值分解。

怎么记忆？

如果大家还是分不清这四个空间，我还有一个小技巧！

大家只需要记住 $XV = US$ 这个式子。

U 和 X 等长，即列向量行数相等，因此 U 一定包含列空间。

U 在矩阵乘积 US 左边，因此包含“左”零空间。

V 和 X 等宽，即行向量列数相等，且 XV 中的 V 是 X 行向量投影方向，因此 V 包含行空间。

V 在矩阵乘积 XV 右边，因此包含“右”零空间。而右零空间，就简称零空间。因为右零空间最常用，所以独占了“零空间”这个更简洁的头衔。

问题来了，要是记不住 $XV = US$ ，怎么办？

就一句话——我们永远 15 岁！

US 代表“我们”， XV 是罗马数字的 15。

23.4 几何视角说空间

下面我们用具体数值从几何视角再强化理解上节介绍的几个空间。

举个例子

给定矩阵 X 如下：

$$X = \begin{bmatrix} 1 & -1 \\ -\sqrt{3} & \sqrt{3} \\ 2 & -2 \end{bmatrix} \quad (12)$$

一眼就能看出来， X 的两个列向量线性相关，因为：

$$x_1 = -x_2 \quad (13)$$

也就是说 X 的秩为 1，即 $\text{rank}(X) = r = 1$ 。

列向量

为了可视化 \mathbf{x}_1 和 \mathbf{x}_2 这两个列向量，我们需要三维直角坐标系 \mathbb{R}^3 ，如图 11 (a) 所示。

白话说， \mathbb{R}^3 才能装下长度为 3 的列向量 \mathbf{x}_1 和 \mathbf{x}_2 。

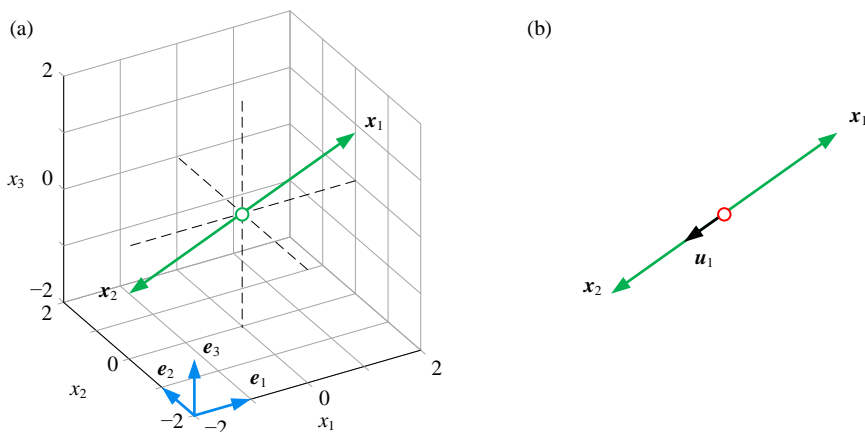


图 11. 从三维空间到一维空间

但是我们发现，实际上，图 11 (b) 告诉我们有了 \mathbf{u}_1 这个单位向量，我们就可以把 \mathbf{x}_1 和 \mathbf{x}_2 写成：

$$\mathbf{x}_1 = a\mathbf{u}_1, \quad \mathbf{x}_2 = b\mathbf{u}_1 \quad (14)$$

也就是说， \mathbb{R}^3 中一维子空间 $\text{span}(\mathbf{u}_1)$ 就足够装下 \mathbf{x}_1 和 \mathbf{x}_2 ，这就是为什么 $\text{rank}(\mathbf{X}) = 1$ 。

那么问题来了，我们如何找到 \mathbf{u}_1 这个单位向量？

根据前文所学，我们知道有至少有两种办法：a) SVD 分解；b) 特征值分解。

SVD 分解

对 \mathbf{X} 进行 SVD 分解得到：

$$\mathbf{X} = \begin{bmatrix} 1 & -1 \\ -\sqrt{3} & \sqrt{3} \\ 2 & -2 \end{bmatrix} = \underbrace{\begin{bmatrix} -0.3536 & -0.9297 & 0.1034 \\ 0.6124 & -0.1465 & 0.7769 \\ -0.7071 & 0.3380 & 0.6211 \end{bmatrix}}_{\mathbf{U}} \underbrace{\begin{bmatrix} 4 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}}_{\mathbf{S}} \underbrace{\begin{bmatrix} -0.7071 & -0.7071 \\ 0.7071 & -0.7071 \end{bmatrix}}_{\mathbf{V}}^T \quad (15)$$

其中，矩阵 \mathbf{U} 的第一列向量就是我们要找的 \mathbf{u}_1 ，而这个 \mathbf{u}_1 便独立张成列空间 $C(\mathbf{X})$ 。

也就是说， $C(\mathbf{X})$ 对应 $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3]$ 线性无关的成分。

顺藤摸瓜，有意思的是完全型 SVD 分解中，我们顺路还得到了 u_2 和 u_3 ，基底 $[u_2, u_3]$ 张起了左零空间 $\text{Null}(X^T)$ 。

而规范正交基 $[u_1, u_2, u_3]$ 则是张成 \mathbb{R}^3 无数规范正交基中的一个。 $[u_1, u_2, u_3]$ 这个独特存在全靠矩阵 X 。

而 $[u_1]$ 和 $[u_2, u_3]$ 补齐得到 \mathbb{R}^3 。显然， u_1 垂直于 u_2 和 u_3 张成的平面 $\text{span}(u_2, u_3)$ 。所示， $[u_1]$ 和 $[u_2, u_3]$ 互为正交补。

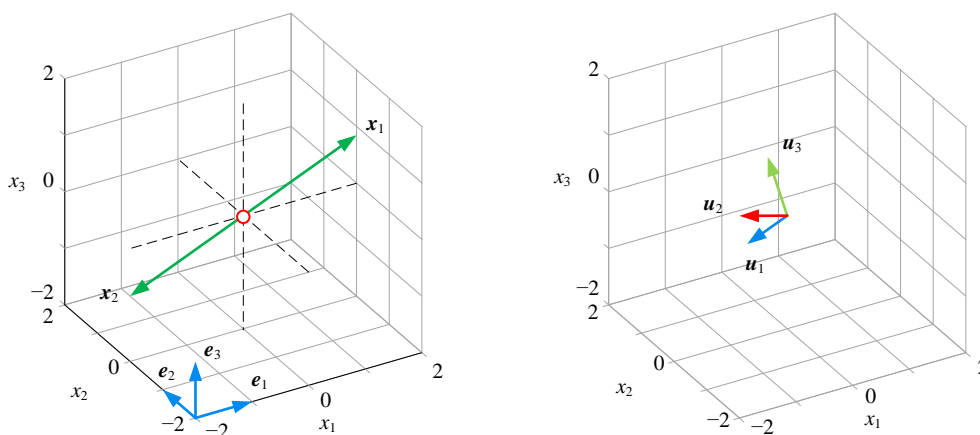


图 12. 矩阵 X 的列空间 $C(X)$ 和左零空间 $\text{Null}(X^T)$

投影

把列向量 x_1 投影到 $U = [u_1, u_2, u_3]$ 中得到：

$$x_1^T U = \begin{bmatrix} 1 & -\sqrt{3} & 2 \end{bmatrix} \begin{bmatrix} \underbrace{\begin{bmatrix} -0.3536 \\ 0.6124 \\ -0.7071 \end{bmatrix}}_{u_1} & \underbrace{\begin{bmatrix} -0.9297 \\ -0.1465 \\ 0.3380 \end{bmatrix}}_{u_2} & \underbrace{\begin{bmatrix} 0.1034 \\ 0.7769 \\ 0.6211 \end{bmatrix}}_{u_3} \end{bmatrix} = \begin{bmatrix} -2.8284 & 0 & 0 \end{bmatrix} \quad (16)$$

也就是说， x_1 在 $[u_1, u_2, u_3]$ 这个标准正交基中的坐标为 $(-2.8282, 0, 0)$ 。

大家可以看到 x_1 在 u_2 和 u_3 上投影结果均为 0，这就是为什么 u_2 和 u_3 上构成左零空间 $\text{Null}(X^T)$ 。

(16) 中的 x_1 转置运算也解释了 $\text{Null}(X^T)$ 括号里面为什么是 X^T 。

同理，把 x_2 投影到 $\{u_1, u_2, u_3\}$ 中得到 x_2 在 $\{u_1, u_2, u_3\}$ 的坐标为 $(2.8282, 0, 0)$ ，对应矩阵运算具体为：

$$\mathbf{x}_2^T \mathbf{U} = \begin{bmatrix} -1 & \sqrt{3} & -2 \end{bmatrix} \begin{bmatrix} \underbrace{\begin{bmatrix} -0.3536 \\ 0.6124 \\ -0.7071 \end{bmatrix}}_{\mathbf{u}_1} & \underbrace{\begin{bmatrix} -0.9297 \\ -0.1465 \\ 0.3380 \end{bmatrix}}_{\mathbf{u}_2} & \underbrace{\begin{bmatrix} 0.1034 \\ 0.7769 \\ 0.6211 \end{bmatrix}}_{\mathbf{u}_3} \end{bmatrix} = \begin{bmatrix} 2.8284 & 0 & 0 \end{bmatrix} \quad (17)$$

特征值分解

当然，我们也可以用特征值分解得到 \mathbf{U} 。首先计算格拉姆矩阵 $\mathbf{X}\mathbf{X}^T$ ：

$$\mathbf{X}\mathbf{X}^T = \begin{bmatrix} 1 & -1 \\ -\sqrt{3} & \sqrt{3} \\ 2 & -2 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -\sqrt{3} & \sqrt{3} \\ 2 & -2 \end{bmatrix}^T = \begin{bmatrix} 2 & -3.4641 & 4 \\ -3.4641 & 6 & -6.9282 \\ 4 & -6.9282 & 8 \end{bmatrix} \quad (18)$$

对 $\mathbf{X}\mathbf{X}^T$ 特征值分解可以得到 \mathbf{U} ：

$$\begin{aligned} \mathbf{X}\mathbf{X}^T &= \begin{bmatrix} 2 & -3.4641 & 4 \\ -3.4641 & 6 & -6.9282 \\ 4 & -6.9282 & 8 \end{bmatrix} \\ &= \underbrace{\begin{bmatrix} -0.3536 & -0.9297 & 0.1034 \\ 0.6124 & -0.1465 & 0.7769 \\ -0.7071 & 0.3380 & 0.6211 \end{bmatrix}}_{\mathbf{U}} \begin{bmatrix} 16 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \underbrace{\begin{bmatrix} -0.3536 & 0.6124 & -0.7071 \\ -0.9297 & -0.1465 & 0.3380 \\ 0.1034 & 0.7769 & 0.6211 \end{bmatrix}}_{\mathbf{U}^T} \end{aligned} \quad (19)$$

图 13 所示为矩阵 \mathbf{X} 的列空间 $C(\mathbf{X})$ 和左零空间 $\text{Null}(\mathbf{X}^T)$ 之间关系。

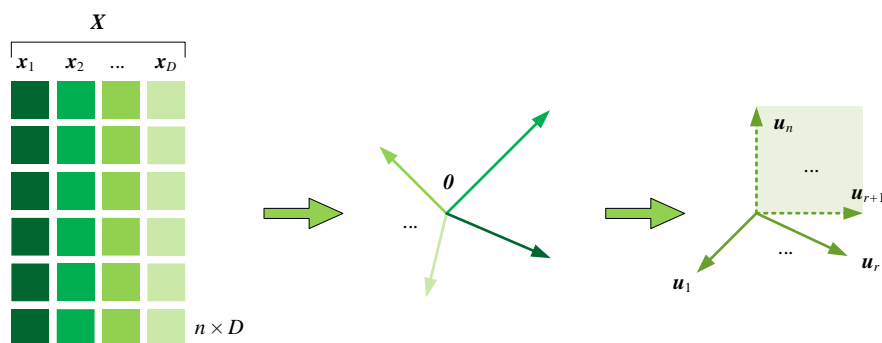


图 13. 矩阵 \mathbf{X} 的列空间 $C(\mathbf{X})$ 和左零空间 $\text{Null}(\mathbf{X}^T)$

行向量

下面，我们聊一下 \mathbf{X} 矩阵的行向量。

很明显 \mathbf{X} 的三个行向量也是线性相关：

$$\mathbf{x}^{(1)} = [1 \ -1], \quad \mathbf{x}^{(2)} = [-\sqrt{3} \ \sqrt{3}], \quad \mathbf{x}^{(3)} = [2 \ -2] \quad (20)$$

如图 14 (a) 所示, 为了装下行向量 $\mathbf{x}^{(1)}$ 、 $\mathbf{x}^{(2)}$ 和 $\mathbf{x}^{(3)}$, 我们需要二维直角坐标系 \mathbb{R}^2 。而图 14 (b) 告诉我们, 用 \mathbf{v}_1 这个单位向量就足以描述 $\mathbf{x}^{(1)}$ 、 $\mathbf{x}^{(2)}$ 和 $\mathbf{x}^{(3)}$, 因为 $\mathbf{x}^{(1)}$ 、 $\mathbf{x}^{(2)}$ 和 $\mathbf{x}^{(3)}$ 可以写成:

$$\mathbf{x}^{(1)} = a\mathbf{v}_1^T, \quad \mathbf{x}^{(2)} = b\mathbf{v}_1^T, \quad \mathbf{x}^{(3)} = c\mathbf{v}_1^T \quad (21)$$

白话说, \mathbb{R}^2 中一维子空间 $\text{span}(\mathbf{v}_1)$ 就足够装下 \mathbf{X} 的三个行向量。

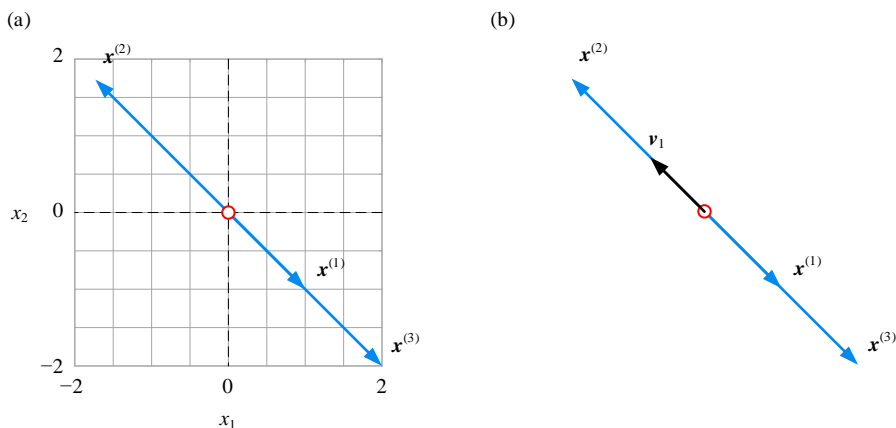


图 14. 从二维空间到一维空间

(15) 给出的 SVD 分解结果已经帮我们找到了 \mathbf{v}_1 。拔出萝卜带出泥, 我们也计算得到 \mathbf{v}_2 。 \mathbf{v}_1 张成行空间 $R(\mathbf{X}) = \text{span}(\mathbf{v}_1)$, \mathbf{v}_2 张成零空间 $\text{Null}(\mathbf{X}) = \text{span}(\mathbf{v}_2)$ 。

而规范正交基 $[\mathbf{v}_1, \mathbf{v}_2]$ 则是张成 \mathbb{R}^2 无数规范正交基中的一个。 $[\mathbf{v}_1, \mathbf{v}_2]$ 是因 \mathbf{X} 而来。

如图 15 (b) 所示, 显然 $R(\mathbf{X}) = \text{span}(\mathbf{v}_1)$ 垂直于 $\text{Null}(\mathbf{X}) = \text{span}(\mathbf{v}_2)$, 即互为正交补。

▲ 格外注意, 大家不要留下错误印象, \mathbf{x}_1 或 $\mathbf{x}^{(1)}$ 就是 \mathbf{u}_1 或 \mathbf{v}_1 的方向重合。一般情况这种重合关系不存在, 本例中产生重合的原因是 $\text{rank}(\mathbf{X}) = 1$ 。

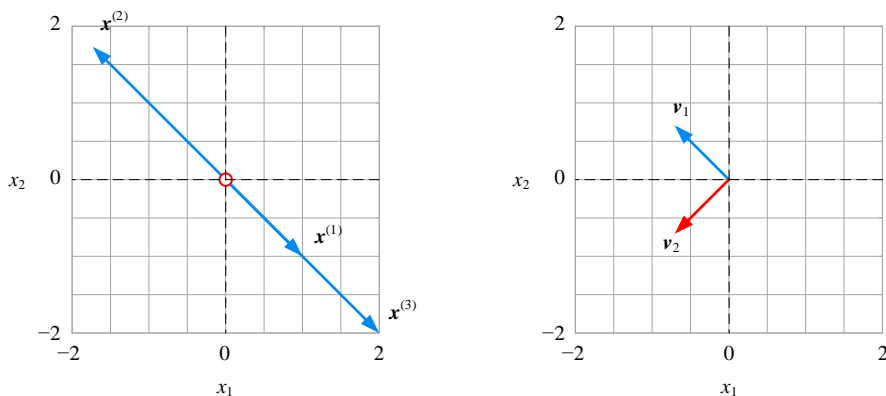


图 15. 矩阵 \mathbf{X} 的行空间 $R(\mathbf{X})$ 和零空间 $\text{Null}(\mathbf{X})$

把行向量 $\mathbf{x}^{(1)}$ 投影到 $[\mathbf{v}_1, \mathbf{v}_2]$ 中得到：

$$\mathbf{x}^{(1)}\mathbf{V} = [1 \quad -1] \left[\underbrace{\begin{bmatrix} -0.7071 \\ 0.7071 \end{bmatrix}}_{\mathbf{v}_1} \underbrace{\begin{bmatrix} -0.7071 \\ -0.7071 \end{bmatrix}}_{\mathbf{v}_2} \right] = [-1.4142 \quad 0] \quad (22)$$

也就是说， $\mathbf{x}^{(1)}$ 在 $[\mathbf{v}_1, \mathbf{v}_2]$ 这个规范正交基中的坐标为 $(-1.4142, 0)$ 。请大家自己计算 $\mathbf{x}^{(2)}$ 和 $\mathbf{x}^{(3)}$ 投影到 $[\mathbf{v}_1, \mathbf{v}_2]$ 结果。

总结来说， $\text{Null}(\mathbf{X})$ 是 \mathbf{X} 的零空间是因为 \mathbf{X} 投影到这个空间的结果都是 0。而 $\text{Null}(\mathbf{X}^T)$ 是 \mathbf{X} 的左零空间是因为， \mathbf{X}^T 投影到这个空间的结果都是 0。

特征值分解

下面，我们再用特征值分解求解 \mathbf{V} 。也是先计算格拉姆矩阵 $\mathbf{X}^T\mathbf{X}$ ：

$$\mathbf{X}^T\mathbf{X} = \begin{bmatrix} 1 & -1 \\ -\sqrt{3} & \sqrt{3} \\ 2 & -2 \end{bmatrix}^T \begin{bmatrix} 1 & -1 \\ -\sqrt{3} & \sqrt{3} \\ 2 & -2 \end{bmatrix} = \begin{bmatrix} 8 & -8 \\ -8 & 8 \end{bmatrix} \quad (23)$$

对 $\mathbf{X}^T\mathbf{X}$ 进行特征值分解，便得到 \mathbf{V} ：

$$\mathbf{X}^T\mathbf{X} = \begin{bmatrix} 8 & -8 \\ -8 & 8 \end{bmatrix} = \underbrace{\begin{bmatrix} -0.7071 & -0.7071 \\ 0.7071 & -0.7071 \end{bmatrix}}_{\mathbf{V}} \underbrace{\begin{bmatrix} 16 & 0 \\ 0 & 0 \end{bmatrix}}_{\Lambda} \underbrace{\begin{bmatrix} -0.7071 & 0.7071 \\ -0.7071 & -0.7071 \end{bmatrix}}_{\mathbf{V}^T} \quad (24)$$

图 16 所示为矩阵 \mathbf{X} 的行空间 $R(\mathbf{X})$ 和零空间 $\text{Null}(\mathbf{X})$ 之间的关系。

此外，值得大家注意的是，比较 (19) 和 (24)，大家容易发现，两个特征值分解都得到了 16 这个特征值。为什么会出现这种情况？下一节将给出答案。

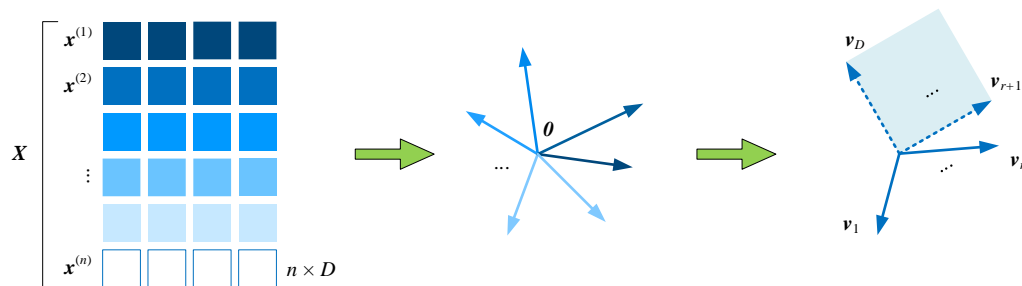


图 16. 矩阵 \mathbf{X} 的行空间 $R(\mathbf{X})$ 和零空间 $\text{Null}(\mathbf{X})$

通过以上分析，希望大家能从几何角度理解六个空间之间的关系。此外，大家也看到奇异值分解的强大之处——任何实数矩阵都可以进行奇异值分解。

23.5 格拉姆矩阵：向量模、夹角余弦值的集合体

我们可以把矩阵 \mathbf{X} 的每一行或每一列分别视作向量。而对于一个向量而言，最能概括它的性质的基本信息莫过于——长度和方向。

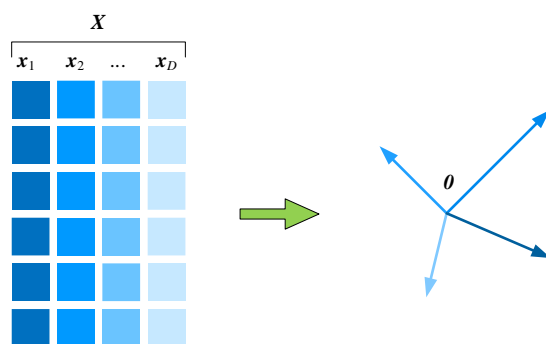


图 17. 矩阵 \mathbf{X} 列向量几何化为空间向量

向量长度不难确定，向量模 (L^2 范数) 就是向量长度。

然而，向量的方向该怎么量化？我们目前接触到几何形体定位最常用手段是平面或三维直角坐标系，直角坐标系在量化位置、长度、方向具有天然优势。

但是对于图 17 所示向量，随着维度不断升高，直角坐标系显得有点力有不逮。

极坐标系

于是，我们想到利用极坐标量化方向。

如图 18 所示，极坐标中定位需要长度和角度，恰巧对应向量的两个重要的元素。唯一的问题是，极坐标系中量化向量和极轴的夹角，即绝对角度。我们接触最多的是向量两两夹角，即相对角度值。

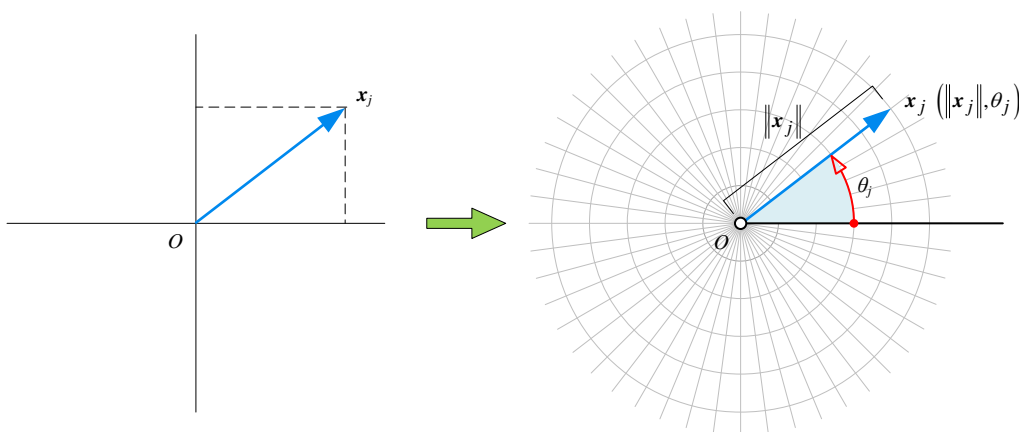


图 18. 从平面直角坐标系到极坐标系，图片参考《数学要素》

此外，向量两两夹角数量也是个问题。数据矩阵 X 有 D 个列向量，这意味着我们可以得到 D 个向量模，以及 $D(D-1)/2$ (C_D^2) 个向量两两夹角余弦值。按照怎样规则保存这些结果？我们反复提到的格拉姆矩阵就是解决方案。而且，本书第 12 章介绍的 Cholesky 分解则帮我们找到这些向量的“绝对位置”。

长度、相对夹角

给定一个 $n \times D$ 数据矩阵 X ，形状细高，也就是 $n > D$ ，它的格拉姆矩阵 G 为：

$$G = X^T X \quad (25)$$

如图 19 所示， G 为对称方阵，形状为 $D \times D$ 。

用向量内积来表达 G ：

$$G = \begin{bmatrix} \langle x_1, x_1 \rangle & \langle x_1, x_2 \rangle & \cdots & \langle x_1, x_D \rangle \\ \langle x_2, x_1 \rangle & \langle x_2, x_2 \rangle & \cdots & \langle x_2, x_D \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle x_D, x_1 \rangle & \langle x_D, x_2 \rangle & \cdots & \langle x_D, x_D \rangle \end{bmatrix} = \begin{bmatrix} \|x_1\| \|x_1\| \cos \theta_{1,1} & \|x_1\| \|x_2\| \cos \theta_{1,2} & \cdots & \|x_1\| \|x_D\| \cos \theta_{1,D} \\ \|x_2\| \|x_1\| \cos \theta_{2,1} & \|x_2\| \|x_2\| \cos \theta_{2,2} & \cdots & \|x_2\| \|x_D\| \cos \theta_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ \|x_D\| \|x_1\| \cos \theta_{D,1} & \|x_D\| \|x_2\| \cos \theta_{D,2} & \cdots & \|x_D\| \|x_D\| \cos \theta_{D,D} \end{bmatrix} \quad (26)$$

可以发现， $G = X^T X$ 包含的信息有两方面： X 列向量的模、列向量两两夹角余弦值。

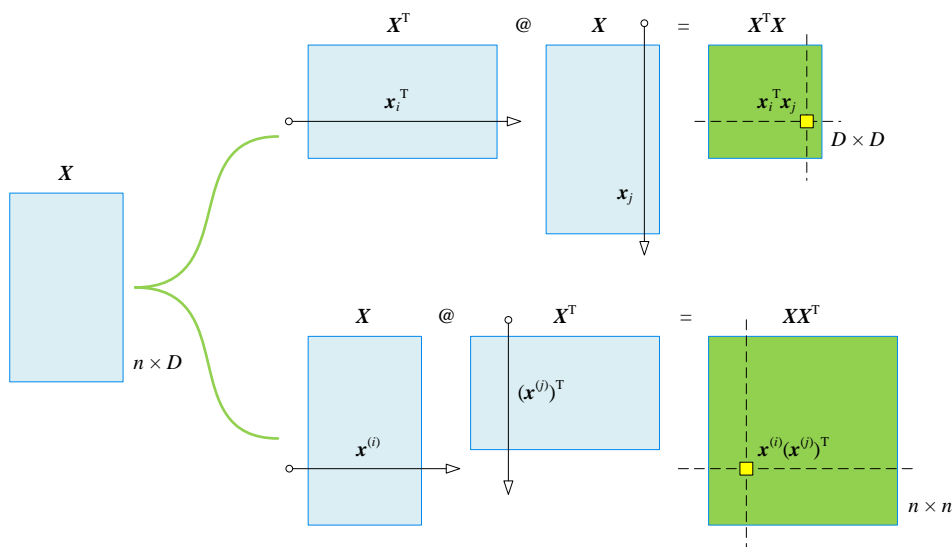


图 19. 两个格拉姆矩阵

而余弦相似度矩阵 C 则进一步减小信息量，只关注列向量夹角余弦值：

$$C = \begin{bmatrix} \frac{\mathbf{x}_1 \cdot \mathbf{x}_1}{\|\mathbf{x}_1\| \|\mathbf{x}_1\|} & \frac{\mathbf{x}_1 \cdot \mathbf{x}_2}{\|\mathbf{x}_1\| \|\mathbf{x}_2\|} & \cdots & \frac{\mathbf{x}_1 \cdot \mathbf{x}_D}{\|\mathbf{x}_1\| \|\mathbf{x}_D\|} \\ \frac{\mathbf{x}_2 \cdot \mathbf{x}_1}{\|\mathbf{x}_2\| \|\mathbf{x}_1\|} & \frac{\mathbf{x}_2 \cdot \mathbf{x}_2}{\|\mathbf{x}_2\| \|\mathbf{x}_2\|} & \cdots & \frac{\mathbf{x}_2 \cdot \mathbf{x}_D}{\|\mathbf{x}_2\| \|\mathbf{x}_D\|} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\mathbf{x}_D \cdot \mathbf{x}_1}{\|\mathbf{x}_D\| \|\mathbf{x}_1\|} & \frac{\mathbf{x}_D \cdot \mathbf{x}_2}{\|\mathbf{x}_D\| \|\mathbf{x}_2\|} & \cdots & \frac{\mathbf{x}_D \cdot \mathbf{x}_D}{\|\mathbf{x}_D\| \|\mathbf{x}_D\|} \end{bmatrix} = \begin{bmatrix} 1 & \cos \theta_{2,1} & \cdots & \cos \theta_{1,D} \\ \cos \theta_{1,2} & 1 & \cdots & \cos \theta_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ \cos \theta_{1,D} & \cos \theta_{2,D} & \cdots & 1 \end{bmatrix} \quad (27)$$

计算 X^T 的格拉姆矩阵，并定义其为 H ：

$$H = XX^T \quad (28)$$

如图 19 所示， H 为对称方阵，形状为 $n \times n$ 。

用向量内积来表达 H ：

$$H = \begin{bmatrix} \langle \mathbf{x}^{(1)}, \mathbf{x}^{(1)} \rangle & \langle \mathbf{x}^{(1)}, \mathbf{x}^{(2)} \rangle & \cdots & \langle \mathbf{x}^{(1)}, \mathbf{x}^{(n)} \rangle \\ \langle \mathbf{x}^{(2)}, \mathbf{x}^{(1)} \rangle & \langle \mathbf{x}^{(2)}, \mathbf{x}^{(2)} \rangle & \cdots & \langle \mathbf{x}^{(2)}, \mathbf{x}^{(n)} \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle \mathbf{x}^{(n)}, \mathbf{x}^{(1)} \rangle & \langle \mathbf{x}^{(n)}, \mathbf{x}^{(2)} \rangle & \cdots & \langle \mathbf{x}^{(n)}, \mathbf{x}^{(n)} \rangle \end{bmatrix} \quad (29)$$

$H = XX^T$ 也包含两方面的信息： X 行向量的模、行向量之间两两夹角余弦值。

特征值分解

下面用特征值分解找到 $X^T X$ 和 XX^T 之间联系。

先对 $G = X^T X$ 进行特征值分解，得到：

$$G = V \Lambda V^T \quad (30)$$

假设 λ_G 为 G 的一个特征值，对应的特征向量为 v ，由此得到等式：

$$Gv = \lambda_G v \quad (31)$$

即，

$$X^T X v = \lambda_G v \quad (32)$$

然后对 H 特征值分解：

$$H = U D U^T \quad (33)$$

U 为特征向量矩阵， D 为特征值对角阵。

假设 λ_H 为 H 的一个特征值，对应特征向量为 u ，构造等式：

$$H u = \lambda_H u \quad (34)$$

即，

$$X X^T u = \lambda_H u \quad (35)$$

(32) 左右乘以 X ，得到：

$$X X^T X v = \lambda_G X v \quad (36)$$

$\underset{u}{\quad} \quad \quad \underset{u}{\quad}$

比较 (35) 和 (36)，可以发现 $X^T X$ 和 XX^T 特征值分解得到的非零特征值存在等价关系。这就回答了为什么 (19) 和 (24) 都有 16 这个特征值这个问题。其实，我们在本书第 16 章也谈过这一现象。

23.6 标准差向量：以数据质心为起点

协方差矩阵可以看成是特殊的格拉姆矩阵，协方差矩阵也是一个“向量模”、“向量间夹角”信息的集合体。

对于形状为 $n \times D$ 的样本数据矩阵 X ，它的协方差矩阵 Σ 可以通过下式计算得到。

$$\Sigma = \frac{\left(\underbrace{X - E(X)}_{\text{Centered}} \right)^T \left(\underbrace{X - E(X)}_{\text{Centered}} \right)}{n-1} = \frac{X_c^T X_c}{n-1} \quad (37)$$

分母上， $n-1$ 仅仅起到取平均作用。 X_c 的格拉姆矩阵为：

$$\mathbf{X}_c^T \mathbf{X}_c = (n-1)\mathbf{\Sigma} \quad (38)$$

图 20 所示, \mathbf{X} 列向量的向量起点为 $\mathbf{0}$ 。而去均值获得 \mathbf{X}_c 过程, 相当于把列向量起点移动到质心 $\mathbf{E}(\mathbf{X})$:

$$\mathbf{X}_c = \mathbf{X} - \mathbf{E}(\mathbf{X}) \quad (39)$$

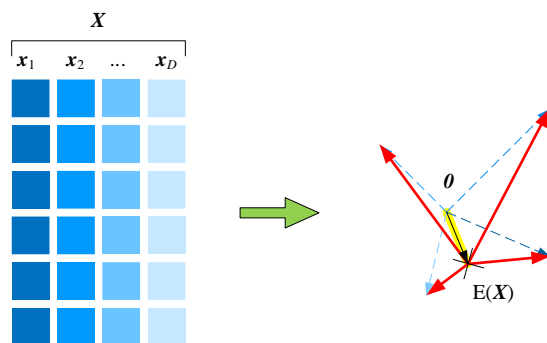


图 20. 数据质心为 \mathbf{X}_c 列向量的起点

将 \mathbf{X}_c 列向量的起点也平移到 $\mathbf{0}$, 和 \mathbf{X} 列向量起点对齐。图 21 比较 \mathbf{X} 和 \mathbf{X}_c 列向量, 显然去均值之后, 向量的长度和向量之间的夹角都发生了变化。有一种特例是, 当质心 $\mathbf{E}(\mathbf{X})$ 本来就在 $\mathbf{0}$ 时, 这样 $\mathbf{X} = \mathbf{X}_c$ 。

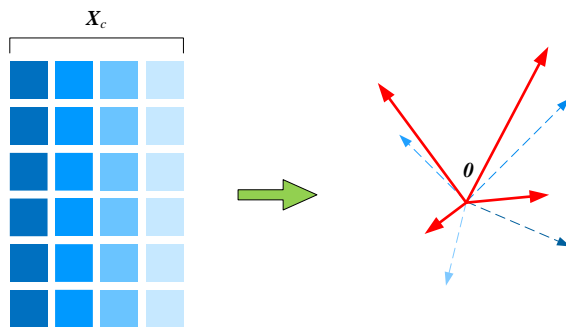


图 21. 比较 \mathbf{X} 和 \mathbf{X}_c 列向量

在数据科学和机器学习应用中, 最常见的三大类数据矩阵就是: 1) 原始数据矩阵 \mathbf{X} ; 2) 中心化数据矩阵 \mathbf{X}_c ; 3) 标准化数据矩阵 \mathbf{Z}_X (z 分数)。

根据本章前文介绍数据矩阵 \mathbf{X} 有四个空间; 显然, 中心化数据矩阵 \mathbf{X}_c 也有自己的四个空间! 那么大家立刻会想到, 标准化数据矩阵 \mathbf{Z}_X , 肯定也有对应的四个空间!

也就是说，如果用 SVD 分解 X 、 X_c 、 Z_X 这三个数据矩阵，会得到不同的结果。下一章则通过各种矩阵分解帮我们分析这三大类数据特点和区别。

标准差向量

整理 (37) 得到 $X_c^T X_c$ ：

$$X_c^T X_c = (n-1)\Sigma = (n-1) \begin{bmatrix} \sigma_1^2 & \rho_{1,2}\sigma_1\sigma_2 & \cdots & \rho_{1,D}\sigma_1\sigma_D \\ \rho_{1,2}\sigma_1\sigma_2 & \sigma_2^2 & \cdots & \rho_{2,D}\sigma_2\sigma_D \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{1,D}\sigma_1\sigma_D & \rho_{2,D}\sigma_2\sigma_D & \cdots & \sigma_D^2 \end{bmatrix} \quad (40)$$

对比 (26) 和 (40)，我们可以把标准差 σ_j 也看做是向量 σ_j ，我们给它起个名字“标准差向量”。

标准差向量 σ_j 之间的夹角的余弦值便是相关性系数。这样 (40) 可以写成：

$$\Sigma = \begin{bmatrix} \langle \sigma_1, \sigma_1 \rangle & \langle \sigma_1, \sigma_2 \rangle & \cdots & \langle \sigma_1, \sigma_D \rangle \\ \langle \sigma_2, \sigma_1 \rangle & \langle \sigma_2, \sigma_2 \rangle & \cdots & \langle \sigma_2, \sigma_D \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle \sigma_D, \sigma_1 \rangle & \langle \sigma_D, \sigma_2 \rangle & \cdots & \langle \sigma_D, \sigma_D \rangle \end{bmatrix} = \begin{bmatrix} \|\sigma_1\| \|\sigma_1\| \cos \phi_{1,1} & \|\sigma_1\| \|\sigma_2\| \cos \phi_{2,1} & \cdots & \|\sigma_1\| \|\sigma_D\| \cos \phi_{1,D} \\ \|\sigma_2\| \|\sigma_1\| \cos \phi_{1,2} & \|\sigma_2\| \|\sigma_2\| \cos \phi_{2,2} & \cdots & \|\sigma_2\| \|\sigma_D\| \cos \phi_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ \|\sigma_D\| \|\sigma_1\| \cos \phi_{1,D} & \|\sigma_D\| \|\sigma_2\| \cos \phi_{2,D} & \cdots & \|\sigma_D\| \|\sigma_D\| \cos \phi_{D,D} \end{bmatrix} \quad (41)$$

如果两个随机变量线性相关，则对应标准差向量平行；如果两个随机变量线性无关，对应的标准差向量正交。

图 22 比较余弦相似度和相关性系数。

▲ 注意，图 22 忽略了 $n-1$ 对缩放的影响。

相关性系数和余弦相似性都描述了两个“相似程度”，也就是靠近的程度；两者取值范围都是 $[-1, 1]$ 。越靠近 1，说明越相似，向量越贴近；越靠近 -1，说明越不同，向量越背离。

不同的是，相关性系数量化“标准差向量” σ_j 之间相似，余弦相似性量化数据矩阵 X 列向量 x_j 之间相似。 x_j 向量的始点为原点 θ ， σ_j 向量始点为数据质心。

大家可能想要知道 x_j 向量和 σ_j 向量到底是什么？它们的具体坐标值又如何？我们下一章回答这个问题。

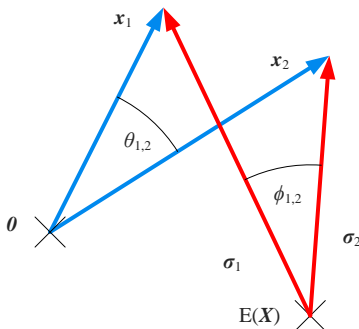


图 22. 余弦相似度和相关性系数的关系，图中忽略标准差向量的缩放系数

相关性系数

类似余弦相似度矩阵 C ，相关性系数矩阵 P 仅仅含有标准差向量夹角（即相关性系数）这一层信息：

$$P = \begin{bmatrix} \frac{\sigma_1 \cdot \sigma_1}{\|\sigma_1\| \|\sigma_1\|} & \frac{\sigma_1 \cdot \sigma_2}{\|\sigma_1\| \|\sigma_2\|} & \cdots & \frac{\sigma_1 \cdot \sigma_D}{\|\sigma_1\| \|\sigma_D\|} \\ \frac{\sigma_2 \cdot \sigma_1}{\|\sigma_2\| \|\sigma_1\|} & \frac{\sigma_2 \cdot \sigma_2}{\|\sigma_2\| \|\sigma_2\|} & \cdots & \frac{\sigma_2 \cdot \sigma_D}{\|\sigma_2\| \|\sigma_D\|} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\sigma_D \cdot \sigma_1}{\|\sigma_D\| \|\sigma_1\|} & \frac{\sigma_D \cdot \sigma_2}{\|\sigma_D\| \|\sigma_2\|} & \cdots & \frac{\sigma_D \cdot \sigma_D}{\|\sigma_D\| \|\sigma_D\|} \end{bmatrix} = \begin{bmatrix} 1 & \cos \phi_{2,1} & \cdots & \cos \phi_{1,D} \\ \cos \phi_{1,2} & 1 & \cdots & \cos \phi_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ \cos \phi_{1,D} & \cos \phi_{2,D} & \cdots & 1 \end{bmatrix} \quad (42)$$

如图 23 所示，以二元随机数为例，相关性系数可以通过散点、二元高斯分布 PDF 曲面、PDF 等高线、椭圆表达。有了本节内容，在众多可视化方案基础上，相关性系数又多了一层几何表达。本系列丛书《概率统计》将讲解随机数、二元高斯分布、概率密度函数 PDF 等概念。此外，《概率统计》中大家会看到无处不在的椭圆。

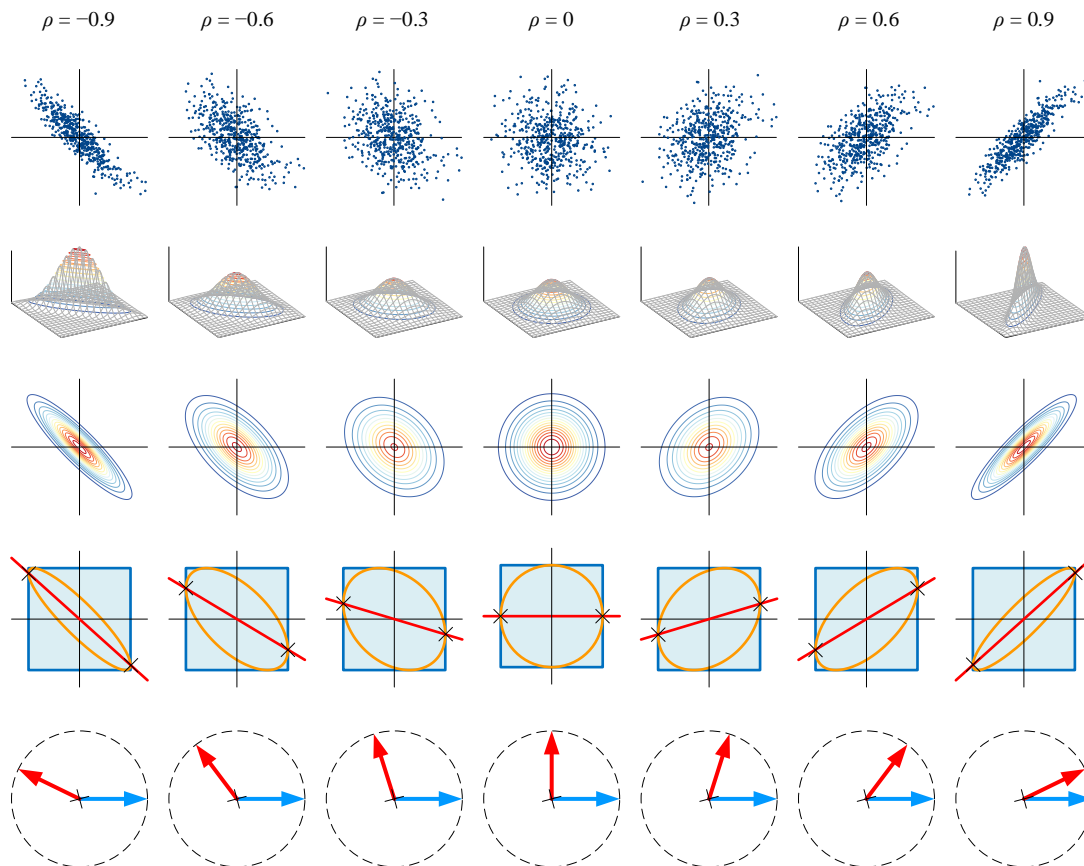


图 23. 相关性系数的几种表达，图中标准差相等，质心位于原点

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

23.7 白话说空间：以鸢尾花数据为例

本章最后一节，我们尝试尽量用大白话把本章之前讲解的四个空间说清楚。本节用的数据是鸢尾花数据前两列，即 $\mathbf{X}_{150 \times 2} = [\mathbf{x}_1, \mathbf{x}_2]$ 。

标准正交基

矩阵 \mathbf{X} 有 150 行、2 列，有 150 个行向量，它们就是图 24 中灰色带箭头的线段。为了装下这 150 个行向量，我们自然而然地想到了 $[\mathbf{e}_1, \mathbf{e}_2]$ —— 平面 \mathbb{R}^2 的标准正交基。

图中散点横坐标就对应 \mathbf{X} 的第一列向量 \mathbf{x}_1 ，纵坐标对应 \mathbf{X} 的第二列向量 \mathbf{x}_2 。

本书第 7 章讲过， $[\mathbf{e}_1, \mathbf{e}_2]$ 表示平面 \mathbb{R}^2 最为自然，因此叫做“标准”正交基。

大家知道，1 维空间是相当于一条过原点的直线，显然图 24 的散点不在一条过原点的直线上。也就是说，要想装下 \mathbf{X} 的行向量至少需要一个二维空间。因此 $[\mathbf{e}_1, \mathbf{e}_2]$ 对于图 24 向量来说，大小正好，没有任何富余。

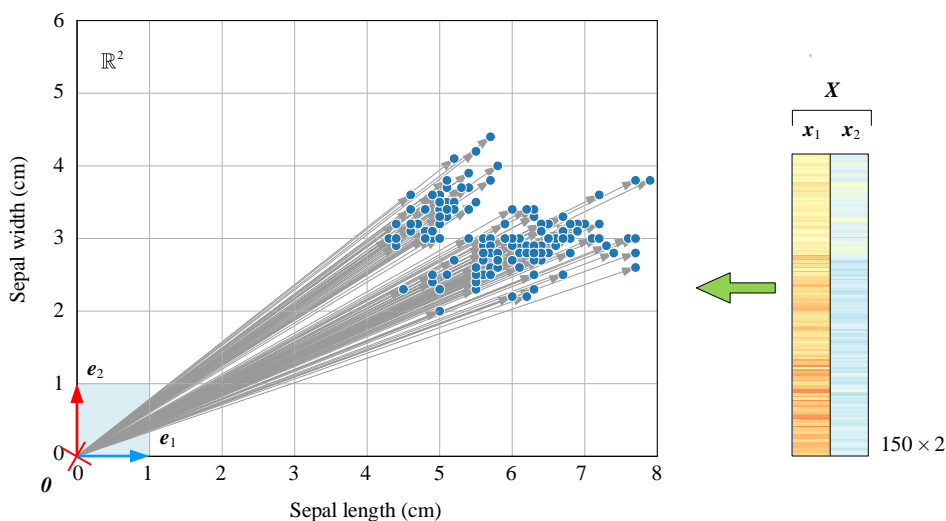
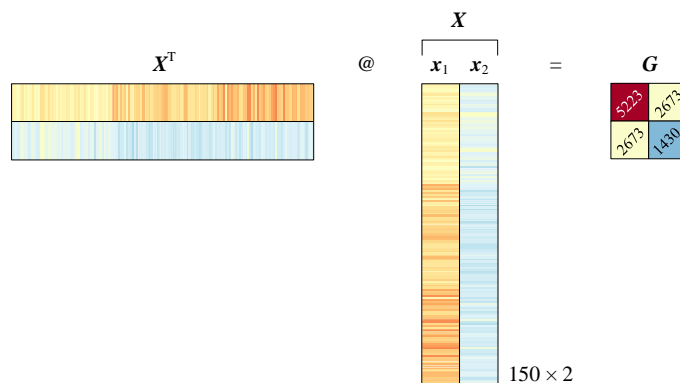


图 24. 找一个能够装下 \mathbf{X} 行向量的空间

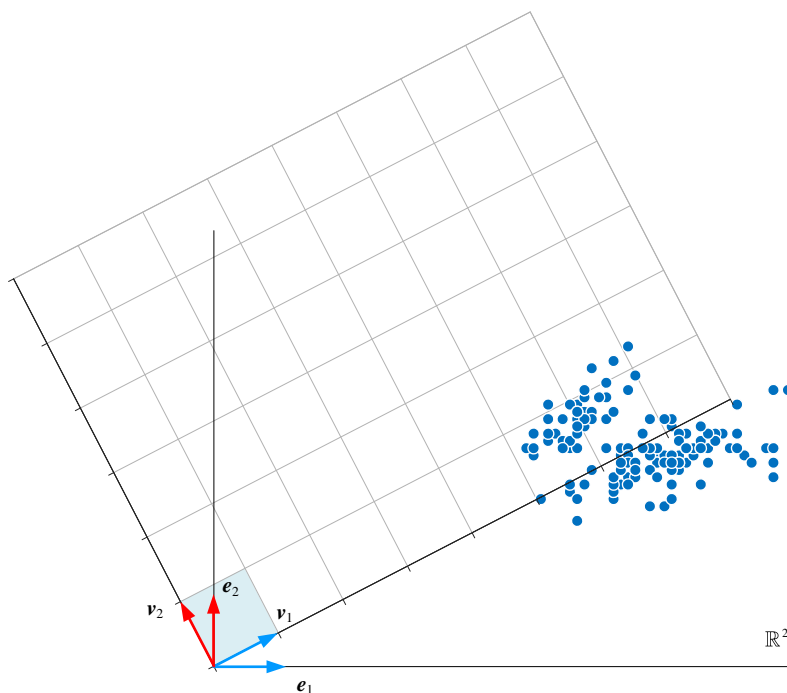
行空间、零空间

根据本章前文所学，为了计算 \mathbf{X} 的**行空间**、**零空间**，我们可以首先计算格拉姆矩阵（如所示），然后对 $\mathbf{X}^T \mathbf{X}$ 特征值分解：

$$\mathbf{G} = \mathbf{X}^T \mathbf{X} = \begin{bmatrix} 5223.85 & 2673.43 \\ 2673.43 & 1430.40 \end{bmatrix} = \underbrace{\begin{bmatrix} 0.888 & -0.459 \\ 0.459 & 0.888 \end{bmatrix}}_{\mathbf{V}} \underbrace{\begin{bmatrix} 6605.05 & & \\ & 49.20 & \\ & & \end{bmatrix}}_{\mathbf{\Lambda}} \underbrace{\begin{bmatrix} 0.888 & 0.459 \\ -0.459 & 0.888 \end{bmatrix}}_{\mathbf{V}^T} \quad (43)$$

图 25. 计算格拉姆矩阵 $\mathbf{X}^T \mathbf{X}$

可以张起 \mathbb{R}^2 的规范正交基有无数个，(43) 中的 $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2]$ 只是其中一个。 $[\mathbf{v}_1, \mathbf{v}_2]$ 在平面上的网格如图 26 所示。 \mathbf{X} 在这个 $[\mathbf{v}_1, \mathbf{v}_2]$ 坐标系中有全新的坐标点。请大家自己回忆怎么计算新的坐标点。

图 26. 规范正交基 $[\mathbf{v}_1, \mathbf{v}_2]$

大家可能会问，之前我们已经在 $[e_1, e_2]$ 这个坐标系中“自然地”描绘了数据矩阵 X ，为何还要劳神费力地寻找 $[v_1, v_2]$ 。

这是因为对 X 来说， $[v_1, v_2]$ 可谓“量身打造”！下面，我们看看 $[v_1, v_2]$ 有何特殊之处。

如图 27 所示， X 向 v_1 投影结果为 $y_1 = Xv_1$ 。 X 是图中的蓝色点 \bullet ， y_1 为图中的蓝色叉 \times 在 $\text{span}(v_1)$ 上的坐标值。蓝色叉 \times 距离原点欧氏距离的平方对应 (43) 中特征值 $\lambda_1 = 6605.05$ 。

利用本书第 18 章介绍的优化视角来观察，给定平面内任意单位向量 v ， $\|Xv\|_2^2$ 的最大值就是 λ_1 。而 $\|Xv\|_2$ 能取得的最大值就是 $\sqrt{\lambda_1}$ ，对应 X 的最大奇异值，即 $s_1 = \sqrt{\lambda_1}$ 。

反之，如图 28 所示， X 向 v_2 投影结果为 $y_2 = Xv_2$ 。给定平面内任意单位向量 v ， $\|Xv\|_2^2$ 的最小值就是 λ_2 。

基底 $[v_1, v_2]$ 对于 X 来说，也显得“捉襟见肘”，维度不能再进一步减小。

如果 X 非列满秩， V 就会出现“余富”，这个余富就是零空间。

比如， $X_1 = Xv_1 \otimes v_1$ 就是图 27 中蓝色叉 \times 在 \mathbb{R}^2 中坐标。蓝色叉 \times 显然都在一条过原点的直线上。 X_1 的秩为 1 也印证了这一点。对于来说， $\text{span}(v_1)$ 足够装下 X_1 ，余富的 $\text{span}(v_2)$ 就是 X_1 的零空间。很明显， X_1 在 $\text{span}(v_2)$ 投影为零向量 θ 。感兴趣的话，大家可以自己计算 X_1 的特征值，它的一个特征值是 $\lambda_1 = 6605.05$ ，另一个特征值为 0。

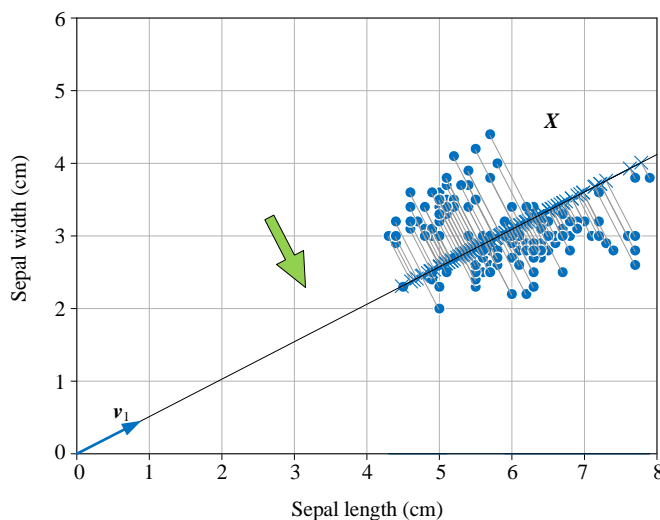


图 27. X 向 v_1 投影

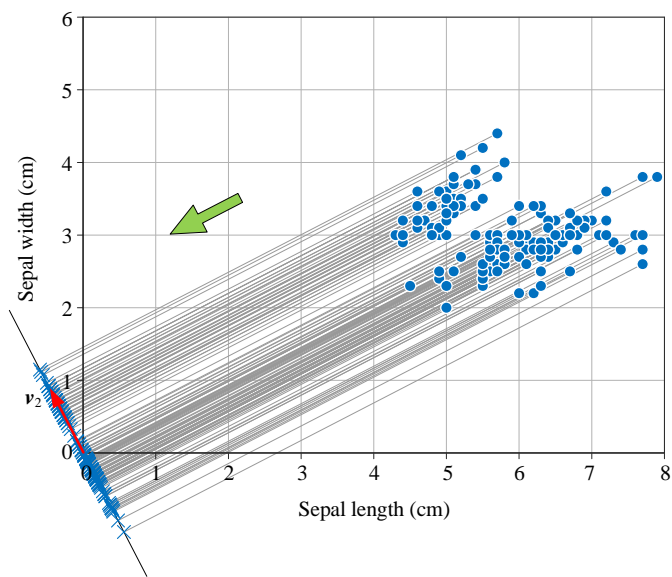


图 28. X 向 v_2 投影

列空间、左零空间

矩阵 X 有两个列向量 x_1 和 x_2 ， x_1 和 x_2 的行数都是 150。为了装下 x_1 和 x_2 ，我们自然想到 \mathbb{R}^{150} 。但是 \mathbb{R}^{150} 对于矩阵 X 来说简直就是“高射炮打蚊子”，小题大做！

下面解释为什么。

为了计算矩阵 X 的列空间、左零空间，我们首先计算格拉姆矩阵 XX^T ，计算过程如图 29 所示。格拉姆矩阵 XX^T 形状为 150×150 。格拉姆矩阵 XX^T 看着很大，实际上它的秩只有 2。也就是说， XX^T 所有 150 个列向量都可以用两个列向量线性组合来表达。

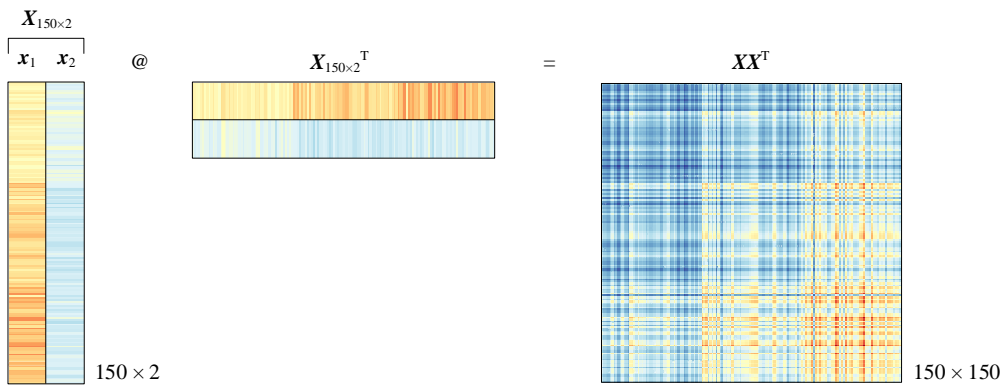


图 29. 计算格拉姆矩阵 XX^T

对 XX^T 特征值分解得到特征向量构成的矩阵 U 如图 30 所示。 U 的形状也是 150×150 。 U 是 \mathbb{R}^{150} 中无数个规范正交阵中的一个。

XX^T 的非零特征值就是 (43) 中的两个特征值，剩余的特征值都为 0。也就是说， U 的前两列 $[u_1, u_2]$ 就是我们要找的列空间， $[u_1, u_2]$ 正好可以装下 X 。剩余的 148 列构成左零空间 $\text{Null}(X^T)$ 。也就是说，想要装下 X 的列向量， \mathbb{R}^{150} 富富有余。

请大家用同样的思路分析 X_c 、 Z_X 。

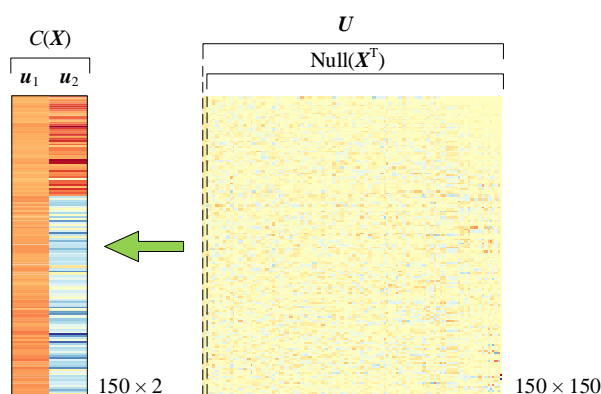


图 30. 格拉姆矩阵 XX^T 的特征向量矩阵 U

有数据的地方，就有矩阵！

有矩阵的地方，就有向量！

有向量的地方，就有几何！

有向量的地方，就有空间！

本书最后三章开启了一场特殊的旅行——“数据三部曲”。这三章梳理总结本书前文核心内容，同时展望这些数学工具的应用。本章作为“数据三部曲”的第一部，主要通过数据矩阵奇异值分解介绍了四个空间。

下图虽然是一幅图，但是其中有四幅子图，它们最能总结本章的核心内容——四个空间。强烈建议大家自行脑补图中缺失的各种符号，以及它们的意义。

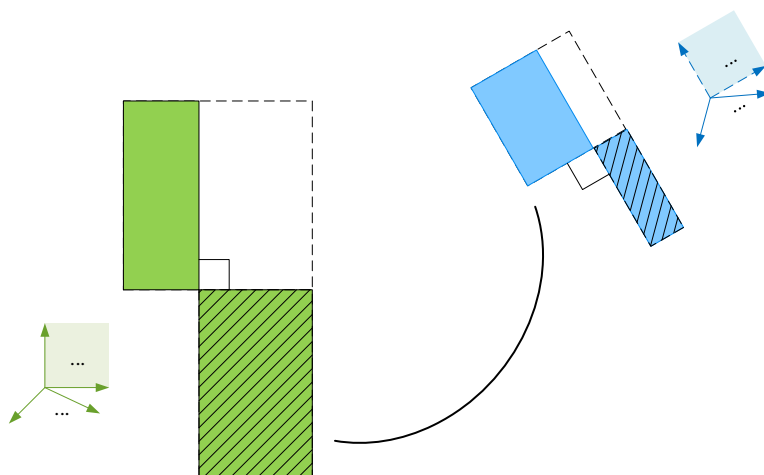


图 31. 总结本章重要内容的四幅图

此外，本书还引出了中心化数据、标准化数据，并创造了“标准差向量”这个概念。格拉姆矩阵是原始数据矩阵列向量长度和两两角度信息的集合体，协方差矩阵则是标准差向量长度和两两角度的结合体。这种类比有助于我们理解线性代数工具在多元统计领域的应用。



推荐大家阅读 MIT 数学教授 Gilbert Strang 的 *Linear Algebra and Learning from Data*。这本书可谓线性代数工具的弹药库，从知识体系上给了本书作者很多启发。图书目前没有免费电子版图书，该书的专属网站提供样张和勘误等资源：

<https://math.mit.edu/~gs/learningfromdata/>