

# 25

## Principal Component Analysis

# 主成分分析

以概率统计、几何、矩阵分解、优化为视角



想象力无边界的人，才创造不可能的事。

*Those who can imagine anything, can create the impossible.*

—— 艾伦·图灵 (Alan Turing) | 英国计算机科学家、数学家，人工智能之父 | 1912 ~ 1954



```
numpy.cov() 计算协方差矩阵
numpy.linalg.eig() 特征值分解
numpy.linalg.svd() 奇异值分解
sklearn.decomposition.PCA() 主成分分析函数
seaborn.heatmap() 绘制热图
seaborn.kdeplot() 绘制 KDE 核概率密度估计曲线
seaborn.pairplot() 绘制成对分析图
numpy.cov() 计算协方差矩阵
numpy.linalg.eig() 特征值分解
numpy.linalg.svd() 奇异值分解
numpy.random.multivariate_normal() 产生多元正态分布随机数
pca.inverse_transform() 将数据 Z 还原成 X
pca.transform(X) 将原始数据转化为数据 Z
seaborn.heatmap() 绘制热图
seaborn.jointplot() 绘制联合分布/散点图和边际分布
seaborn.kdeplot() 绘制 KDE 核概率密度估计曲线
seaborn.pairplot() 绘制成对分析图
sklearn.decomposition.PCA() 主成分分析函数
```



## 25.1 再聊主成分分析

主成分分析 (Principal Component Analysis, PCA) 是重要的降维工具。PCA 可以显著减少数据的维数，同时保留数据中对方差贡献最大的成分。另外对于多维数据，PCA 可以作为一种数据可视化的工具。PCA 还可以用来构造回归模型，这是《数据有道》一册要介绍的内容。

本章将以概率统计、几何、矩阵分解、优化为视角给大家全景展示主成分分析。此外，这一章大家可以把它看成丛书“数学”板块的一个总结。

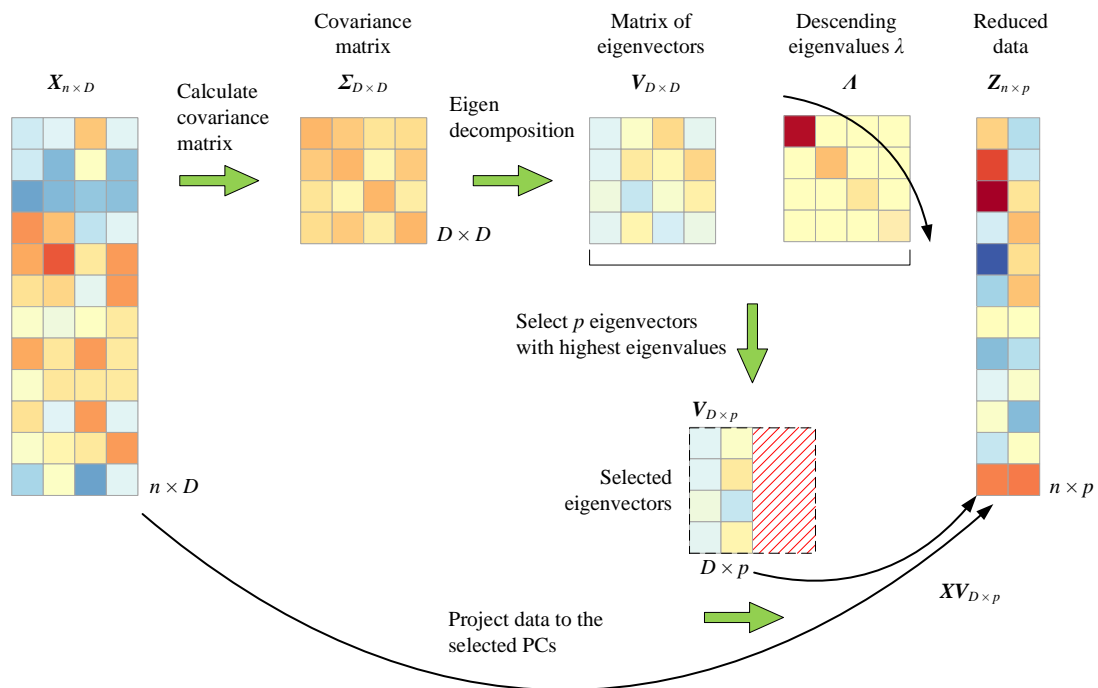


图 1. 主成分分析过程；技术路线：特征值分解协方差矩阵

PCA 的一般步骤如下：

- ▶ 对原始数据  $X_{n \times D}$  的协方差矩阵  $\Sigma_{D \times D}$ ;
- ▶ 计算  $\Sigma$  特征值  $\lambda_i$  与特征向量矩阵  $V_{D \times D}$ ，即因子载荷；
- ▶ 对特征值  $\lambda_i$  从大到小排序，选择其中特征值最大的  $p$  个特征向量；
- ▶ 将原始数据投影到这  $p$  个正交向量构建的新空间中，得到因子得分  $Z_{n \times p}$ 。

很多时候，在第一步中，我们直接对原始数据进行标准化 (standardization) 处理，即计算  $X$  的  $z$  分数。标准化防止不同特征上方差差异过大。我们在《矩阵力量》第 25 章看到的的就是利用标准化数据进行 PCA 分析的技术路线。标准化数据的协方差矩阵实际上就是原数据的相关性系数矩阵。

而有些情况，对原始数据  $\mathbf{X}_{n \times D}$  进行中心化(去均值)就足够了，即将数据质心移到原点。

图 1 所示为通过分解协方差矩阵进行主成分分析过程；当然，也可以通过奇异值分解中心化数据  $\mathbf{X}_c$  进行主成分分析。本章将主要介绍通过特征值分解协方差矩阵进行主成分分析。

## 25.2 原始数据

《矩阵力量》介绍过，样本数据矩阵  $\mathbf{X}$  可以分别通过行和列来解释。矩阵  $\mathbf{X}$  每一列代表一个特征向量：

$$\mathbf{X} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \mathbf{x}_3 \quad \mathbf{x}_4] \quad (1)$$

$\mathbf{X}$  矩阵每一行代表一个样本。比如， $\mathbf{X}$  矩阵第一行对应是第一个数据点，它写成一个行向量  $\mathbf{x}^{(1)}$ ：

$$\mathbf{x}^{(1)} = [x_{1,1} \quad x_{1,2} \quad x_{1,3} \quad x_{1,4}] \quad (2)$$

图 2 展示原始数据矩阵  $\mathbf{X}$  热图，红色色系代表正数，蓝色色系代表负数，黄色接近 0。 $\mathbf{X}$  矩阵有 12 行，即 12 个样本； $\mathbf{X}$  矩阵有 4 列，即 4 个特征。

注意，本例中假设  $\mathbf{X}$  已经中心化  $E(\mathbf{X}) = 0$ ，即质心位于原点。

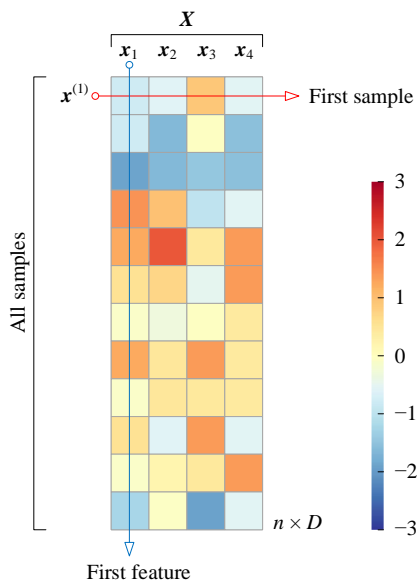


图 2. 原始数据  $\mathbf{X}$  热图， $D = 4$ ， $n = 12$ ， $\mathbf{X}$  已经去均值

### 分布特征

图 3 所示为矩阵  $X$  每一列特征数据的分布情况；可以发现它们之间的均方差区别不大。但是经过主成分分解之后，大家可以明显发现每一列新特征数据均方差大小完全不同。

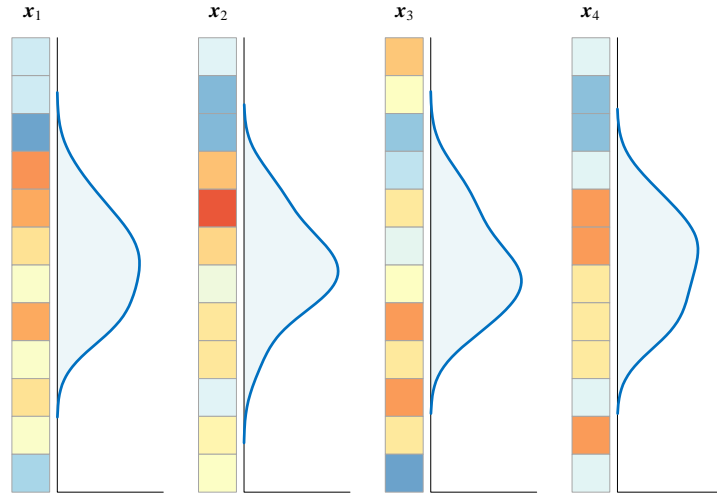


图 3.  $X$  四个特征向量数据分布

## 25.3 特征值分解协方差矩阵

本书第 13 章介绍过， $X$  的协方差矩阵  $\Sigma$  可以通过下式计算得到：

$$\Sigma = \frac{(X - E(X))^T (X - E(X))}{n-1} = \frac{X_c^T X_c}{n-1} \quad (3)$$

其中， $E(X)$  也常被称作原始数据  $X$  的质心； $X - E(X)$  相当于数据中心化。当  $n$  足够大，(3) 的分母可以用  $n$  替换。本例设定  $E(X) = \mathbf{0}^T$ ，即  $X = X_c$ 。

如图 5 所示， $\Sigma$  为实数对称矩阵，它的特征值分解（谱分解）可以写作：

$$\Sigma = V \Lambda V^T \quad (4)$$

$V$  为正交矩阵。 $V$  和自己转置  $V^T$  乘积为单位阵  $I$ ，即：

$$V^T V = I \quad (5)$$

特征值方阵  $\Lambda$  主对角线元素为特征值  $\lambda$ ，特征值从大到小排列：

$$\Lambda = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_d \end{bmatrix}, \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \quad (6)$$

本书前文介绍过，从统计学角度来讲， $\lambda_j$  是第  $j$  个主成分所贡献的方差。

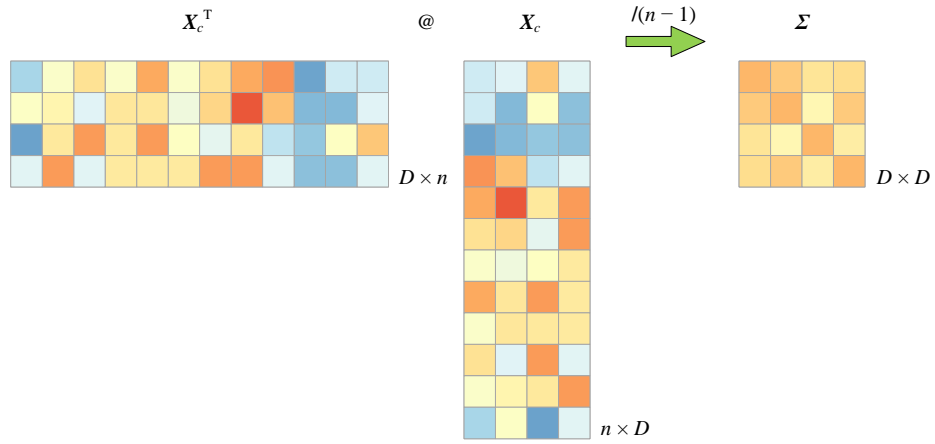


图 4. 计算原始数据协方差矩阵,  $D=4$ ,  $n=12$

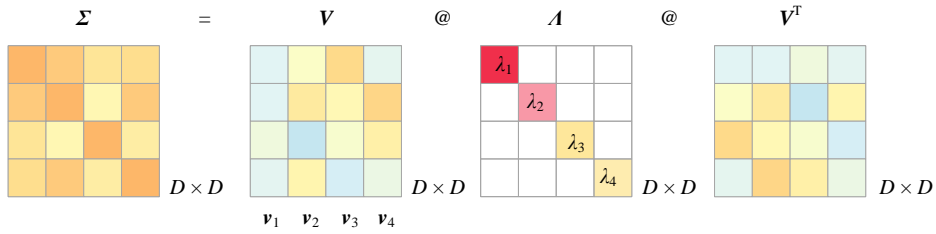


图 5. 协方差矩阵特征值分解,  $D=4$

## 主成分、因子载荷

$V$  为特征向量构造的  $D \times D$  的方阵:

$$V = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_D \\ \text{PC1} & \text{PC2} & & \end{bmatrix} = \begin{bmatrix} v_{1,1} & v_{1,2} & \cdots & v_{1,D} \\ v_{2,1} & v_{2,2} & \cdots & v_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ v_{D,1} & v_{D,2} & \cdots & v_{D,D} \end{bmatrix} \quad (7)$$

$\mathbf{v}_1$  被称作第一主成分 (first principal component), 本书常记做 PC1;  $\mathbf{v}_2$  被称作第二主成分 (second principal component), 记做 PC2; 以此类推。

$V$  的列向量也叫因子载荷 (loadings)。注意, 有些文献中因子载荷为:

$$V\sqrt{\Lambda} = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \cdots \quad \mathbf{v}_D] \begin{bmatrix} \sqrt{\lambda_1} & & & \\ & \sqrt{\lambda_2} & & \\ & & \ddots & \\ & & & \sqrt{\lambda_D} \end{bmatrix} = [\sqrt{\lambda_1}\mathbf{v}_1 \quad \sqrt{\lambda_2}\mathbf{v}_2 \quad \cdots \quad \sqrt{\lambda_D}\mathbf{v}_D] \quad (8)$$

## 迹, 总方差

本书前文介绍过，协方差矩阵  $\Sigma$  的迹  $\text{trace}(\Sigma)$  等于的特征值方阵  $\Lambda$  迹  $\text{trace}(\Lambda)$ ：

$$\text{trace}(\Sigma) = \sigma_1^2 + \sigma_2^2 + \cdots + \sigma_D^2 = \sum_{j=1}^D \sigma_j^2 = \text{trace}(\Lambda) = \lambda_1 + \lambda_2 + \cdots + \lambda_D = \sum_{j=1}^D \lambda_j \quad (9)$$

第  $j$  个特征值  $\lambda_j$  对方差总和 (total variance) 的贡献百分比为：

$$\frac{\lambda_j}{\sum_{i=1}^D \lambda_i} \times 100\% \quad (10)$$

前  $p$  个特征值，即  $p$  个主成了解释总方差 (total variance explained) 的百分比为：

$$\frac{\sum_{j=1}^p \lambda_j}{\sum_{i=1}^D \lambda_i} \times 100\% \quad (11)$$

主成分分析中，我们常用陡坡图 (scree plot) 可视化这个百分比，《数据有道》一册中大家会看到很多实例。

## 25.4 投影

本节从投影角度介绍 PCA。数据矩阵  $X$  投影到矩阵  $V$  正交系 ( $v_1, v_2, \dots, v_D$ ) 得到新特征数据矩阵  $Z$ ，即：

$$Z = XV \quad (12)$$

$V$  常被称作**因子载荷** (factor loadings)， $Z$  常被称作**因子得分** (factor score)。图 6 所示  $Z = XV$  矩阵运算原理图。《矩阵力量》第 10 章特别介绍过这种数据投影，建议大家回顾。

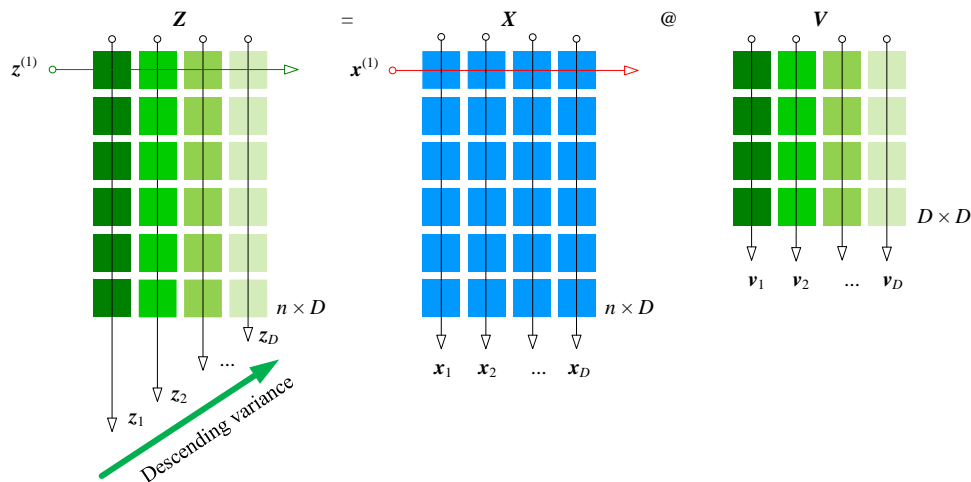


图 6. PCA 分解数据关系  $Z = XV$

图 7 所示为将图 2 给出数据矩阵  $X$  投影到矩阵  $V$ ，得到新特征数据矩阵  $Z$ 。

值得强调的一点是，如果  $X$  没有中心化，把原始数据  $X$  或中心化数据  $X_c$  投影到  $V$  中结果不一样。从统计角度来看，差异主要体现在质心位置，而投影得到的数据协方差矩阵相同。

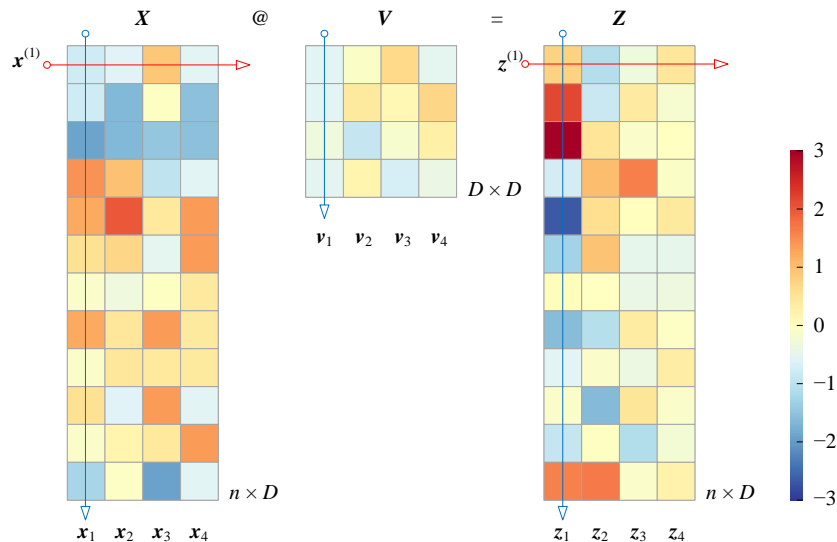


图 7.  $Z$ 、 $X$  和  $V$  这三个矩阵关系和热图

## $Z$ 的列向量

前文讨论过，矩阵  $X$  每一列特征数据方差区别不大 (见图 3)；而图 8 告诉我们，经过 PCA 分解得到的矩阵  $Z$  四个新特征数据分布差异巨大。

如图 8 所示，第一列  $z_1$  数据分布最为分散，也就是第一主成分 (first principal component) 解释了数据中最多方差。第一列  $z_1$  到第四列  $z_4$  数据分散情况逐渐降低，热图对应的色差从明显到模糊。

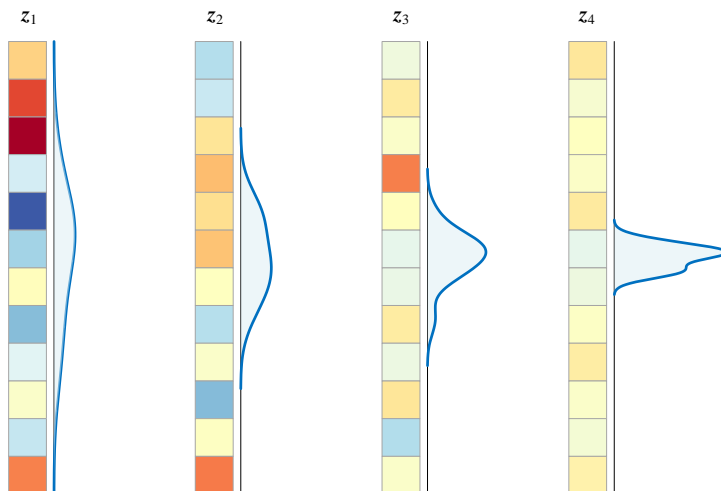


图 8.  $\mathbf{Z}$  四个新特征数据分布

将 (12) 展开得到：

$$\begin{bmatrix} z_1 & z_2 & \cdots & z_D \end{bmatrix} = \mathbf{X} \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_D \\ \text{PC1} & \text{PC2} & & \end{bmatrix} \quad (13)$$

由此，得到图 9 所示主成分分析运算的数据关系：

$$\begin{cases} z_1 = \mathbf{X}\mathbf{v}_1 \\ z_2 = \mathbf{X}\mathbf{v}_2 \\ \vdots \\ z_D = \mathbf{X}\mathbf{v}_D \end{cases} \quad (14)$$

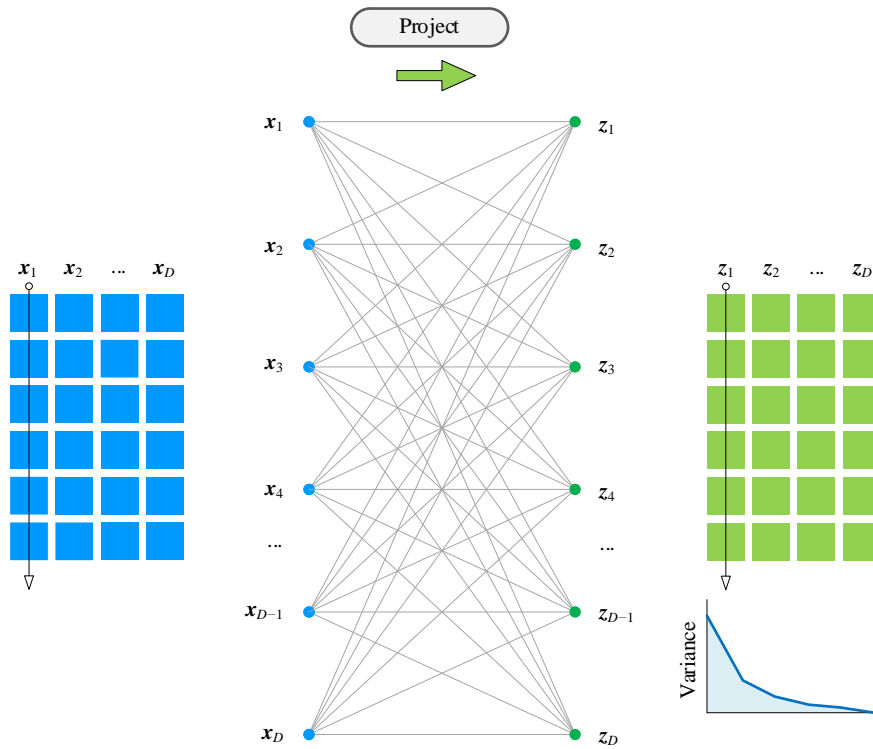


图 9. PCA 运算数据关系

注意， $z_i$  和  $z_j$  ( $i \neq j$ ) 相互正交：

$$\begin{aligned} z_i^T z_j &= (\mathbf{X}\mathbf{v}_i)^T \mathbf{X}\mathbf{v}_j = \mathbf{v}_i^T \mathbf{X}^T \mathbf{X}\mathbf{v}_j = (n-1)\mathbf{v}_i^T \boldsymbol{\Sigma} \mathbf{v}_j \\ &= (n-1)\mathbf{v}_i^T \boldsymbol{\Sigma} \mathbf{v}_j = (n-1)\mathbf{v}_i^T \mathbf{V} \boldsymbol{\Lambda} \mathbf{V}^T \mathbf{v}_j = 0 \end{aligned} \quad (15)$$

## 线性组合



如图 10 所示，列向量  $\mathbf{v}_1$  每一个元素相当于原始数据  $\mathbf{X}$  每一列  $[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_D]$  对应的系数。

以第一主成分数据  $z_1$  为例，将  $\mathbf{X}$  向  $\mathbf{v}_1$  投影，展开得到：

$$\mathbf{z}_1 = \mathbf{X}\mathbf{v}_1 \quad (16)$$

(16) 展开得到：

$$\mathbf{z}_1 = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_D] \begin{bmatrix} v_{1,1} \\ v_{2,1} \\ \vdots \\ v_{D,1} \end{bmatrix} = v_{1,1}\mathbf{x}_1 + v_{2,1}\mathbf{x}_2 + \dots + v_{D,1}\mathbf{x}_D \quad (17)$$

$\mathbf{v}_1, \text{PC1}$

白话讲， $z_1$  相当于  $[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_D]$  每一列的混合成分，即线性组合。

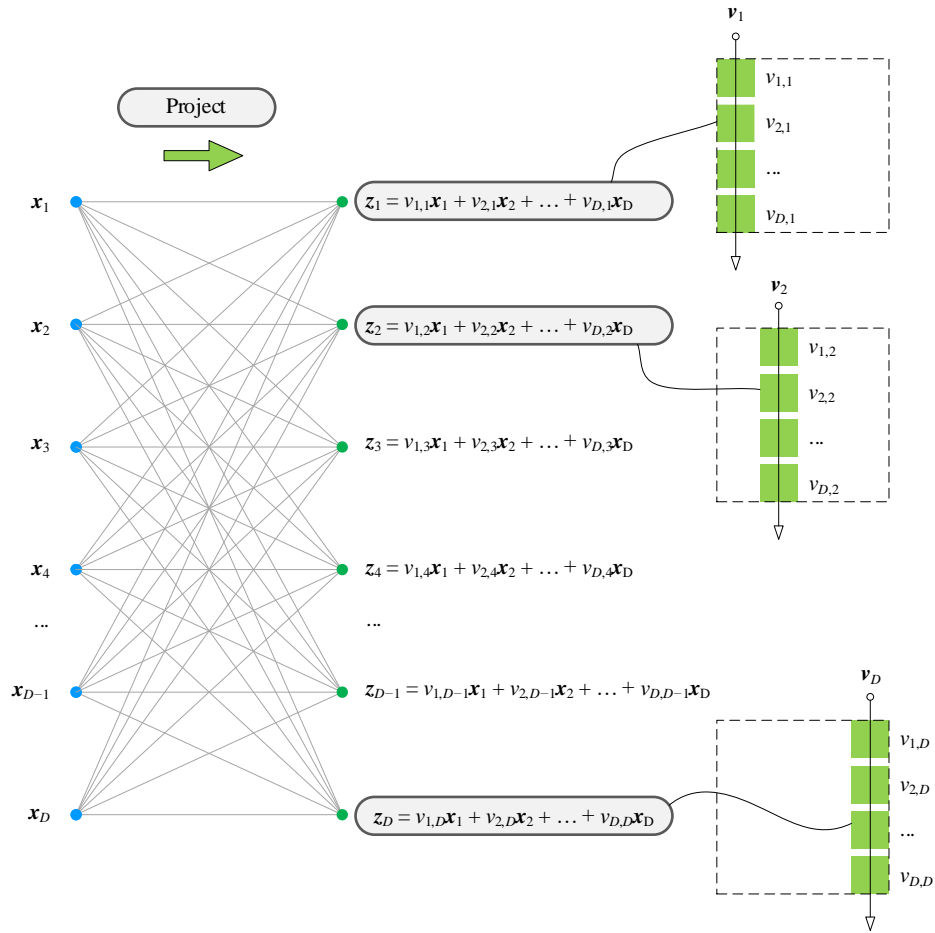


图 10.  $\mathbf{V}$  每一列元素相当于原始数据  $\mathbf{X}$  每一列系数

## 朝向量投影

图 11 所示  $z_1 = Xv_1$  运算相当于数据  $X$  向  $v_1$  向量 (第一主成分) 投影获得  $z_1$ , 即一个四特征数据  $X$  投影到  $v_1$  得到一维新特征数据。图 12 展示  $z_2 = Xv_2$  运算等价于数据  $X$  向  $v_2$  (第二主成分) 投影获得  $z_2$ 。图 11 ~ 图 14 分别展示数据矩阵  $X$  向  $v_1$ 、 $v_2$ 、 $v_3$  和  $v_4$  向量投影。

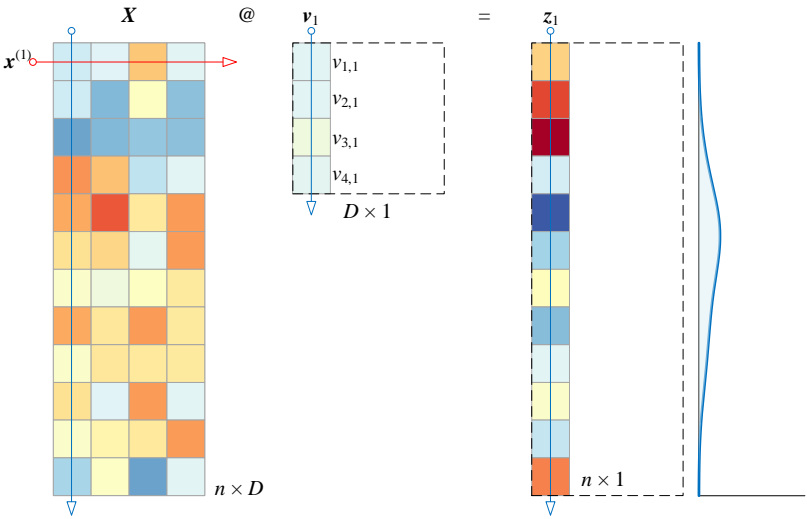


图 11. 数据  $X$  向  $v_1$  向量投影

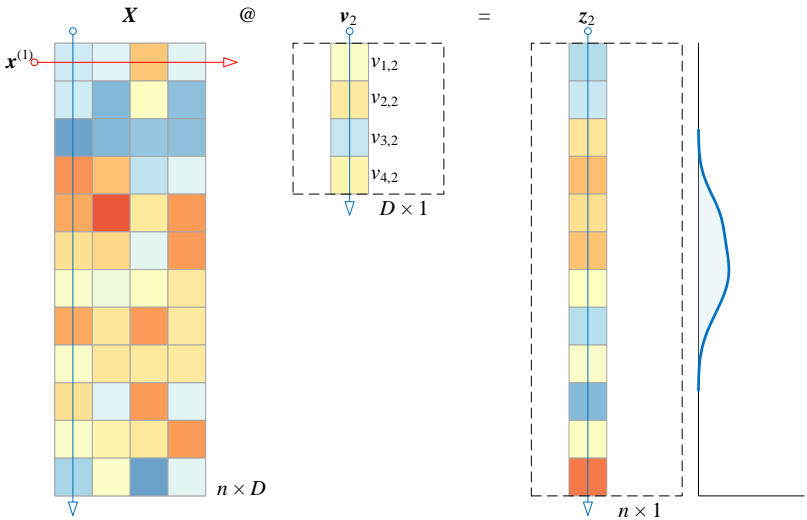


图 12. 数据  $X$  向  $v_2$  向量投影

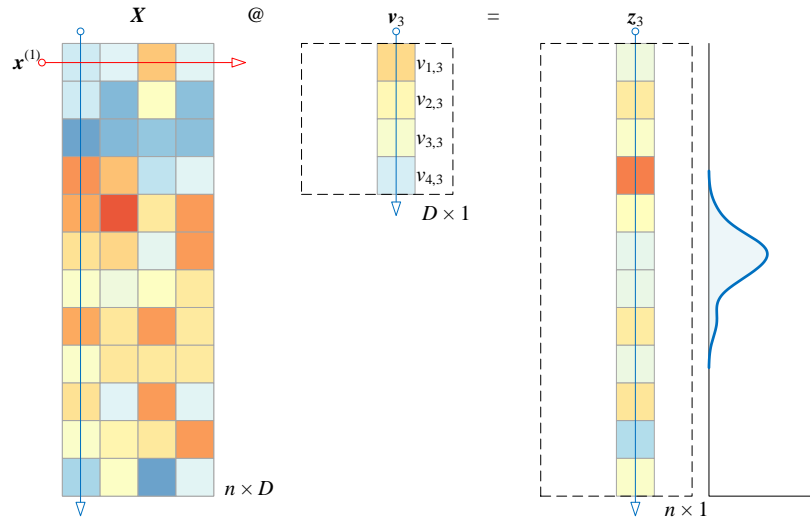


图 13. 数据  $X$  向  $v_3$  向量投影

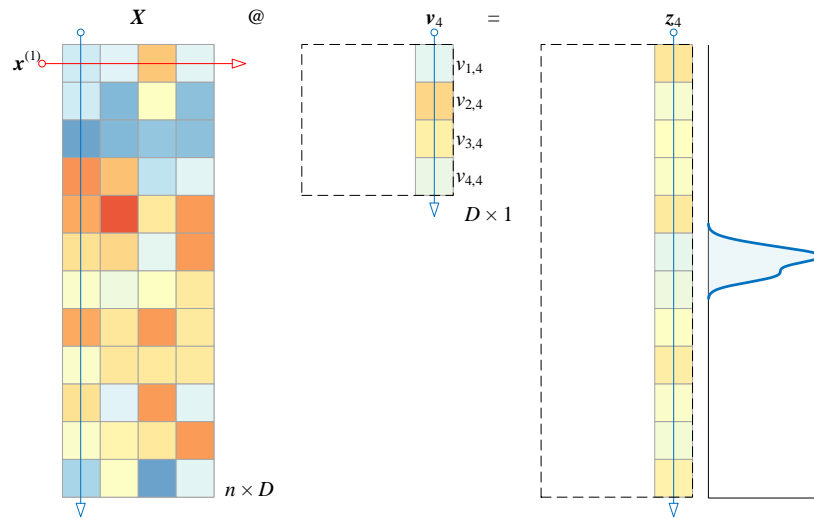


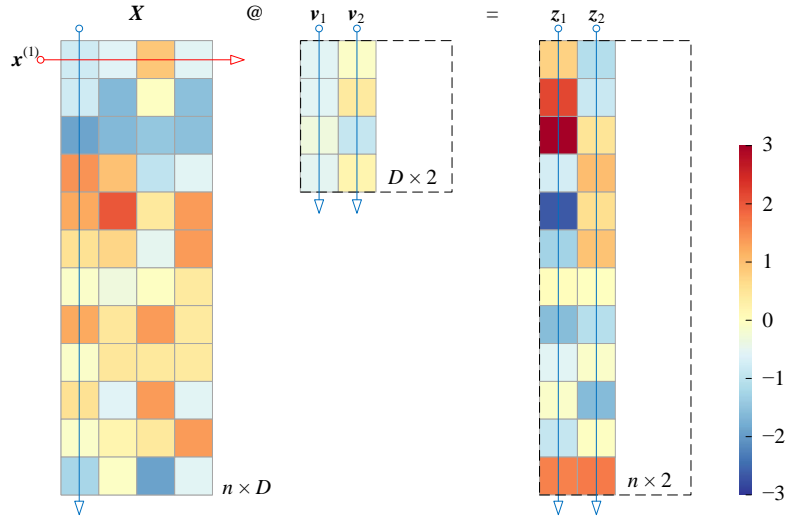
图 14. 数据  $X$  向  $v_4$  向量投影

## 朝平面投影

同样， $[z_1, z_2]$  是数据  $X$  向  $[v_1, v_2]$  投影结果，即相当于四维数据  $X$  向正交平面投影。运算过程如下：

$$\begin{bmatrix} z_1 & z_2 \end{bmatrix} = X \begin{bmatrix} v_1 & v_2 \end{bmatrix} \quad (18)$$

图 15 所示为 (18) 运算过程及结果热图。

图 15. 数据  $X$  向  $[v_1, v_2]$  投影

### $Z$ 的协方差矩阵

前文假设  $X$  已经中心化，因此  $z_1$  的期望值为 0。对  $z_1$  求方差，可以得到：

$$\text{var}(z_1) = \frac{(Xv_1)^T (Xv_1)}{n-1} = \frac{v_1^T X^T X v_1}{n-1} = v_1^T \underbrace{\frac{X^T X}{n-1}}_{\Sigma} v_1 = v_1^T \Sigma v_1 \quad (19)$$

类似地，

$$\text{var}(z_2) = v_2^T \Sigma v_2, \quad \dots, \quad \text{var}(z_D) = v_D^T \Sigma v_D \quad (20)$$

这样，整个新特征数据  $Z$  的协方差矩阵，可以通过下式计算得到：

$$\begin{aligned} \text{var}(Z) &= \frac{(XV)^T (XV)}{n-1} = \frac{V^T X^T X V}{n-1} \\ &= V^T \underbrace{\frac{X^T X}{n-1}}_{\Sigma} V = V^T \Sigma V = \begin{bmatrix} v_1^T \Sigma v_1 & & \\ & v_2^T \Sigma v_2 & \\ & & \ddots \\ & & & v_D^T \Sigma v_D \end{bmatrix} = \Lambda = \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_D \end{bmatrix} \end{aligned} \quad (21)$$

观察 (21) 所示协方差矩阵，可以发现主对角线以外元素均为 0，也就是  $Z$  的列向量两两正交，线性相关系数为 0。

$Z_{n \times p}$  的协方差矩阵为：

$$\text{var}(Z_{n \times p}) = \frac{(XV_{D \times p})^T (XV_{D \times p})}{n-1} = V_{D \times p}^T \frac{X^T X}{n-1} V_{D \times p} = V_{D \times p}^T \Sigma V_{D \times p} = \Lambda_{p \times p} = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{bmatrix} \quad (22)$$

对于投影数据的方差计算，我们已经在本书第 14 章详细介绍过，请感兴趣的读者自行回顾复习。

## 25.5 几何视角看 PCA

如图 16 所示，椭圆中心对应质心  $\mu$ ，椭圆和  $\pm\sigma$  标准差构成的正方形相切，四个切点分别为  $A$ 、 $B$ 、 $C$  和  $D$ ，对角切点两两相连得到两条直线  $AC$ 、 $BD$ 。

$AC$  相当于在给定  $X_2$  条件下  $X_1$  的条件概率期望值； $BD$  相当于在给定  $X_1$  条件下  $X_2$  的条件概率期望值。

图 16 中， $EF$  为椭圆长轴； $FH$  为椭圆短轴。而  $EF$  就相当于本章介绍的第一主成分， $FH$  为第二主成分。

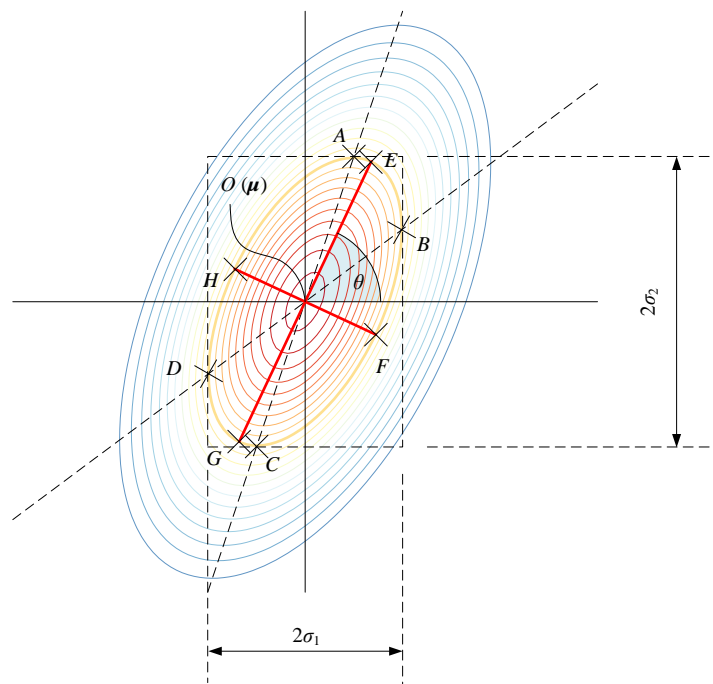


图 16. 椭圆和  $\pm\sigma$  标准差长方形的关系

$\Sigma$  特征值分解进行主成分分析的具体步骤如图 17 所示。假设图 17 原始数据已经标准化，计算得到协方差矩阵  $\Sigma$ ，找到  $\Sigma$  对应椭圆的半长轴所在方向  $v_1$ 。 $v_1$  对应的便是第一主成分 (first principal component)。原始数据朝  $v_1$  投影得到的数据对应最大方差。

整个过程实际上用到了我们在丛书《矩阵力量》一本中介绍的平移、缩放、正交化、投影、旋转等线性变换操作。

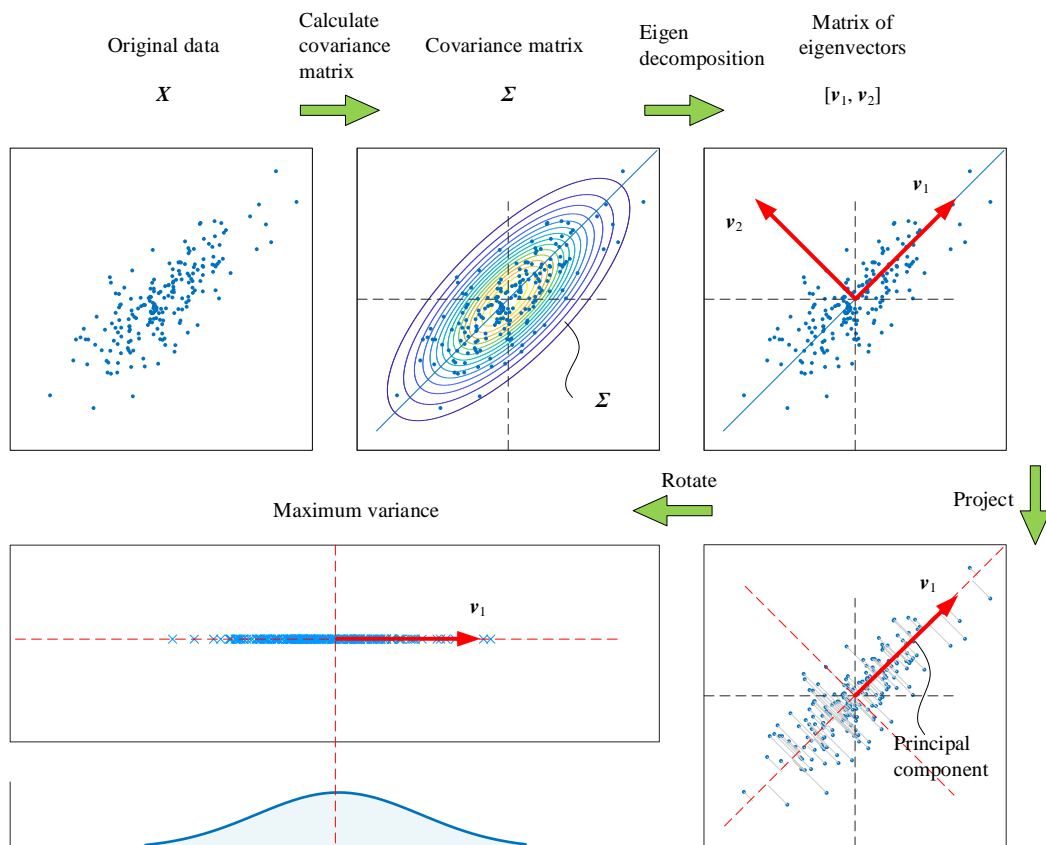


图 17. 几何视角下通过特征值分解协方差矩阵进行主成分分析

从线性变换角度来看，主成分分析无非就是，在不同的坐标系中看同一组数据。如图 18 所示，数据朝不同方向投影会得到不同的投影结果，对应不同的分布；朝椭圆长轴方向投影，得到的数据均方差最大；朝椭圆短轴方向投影得到的数据均方差最小。

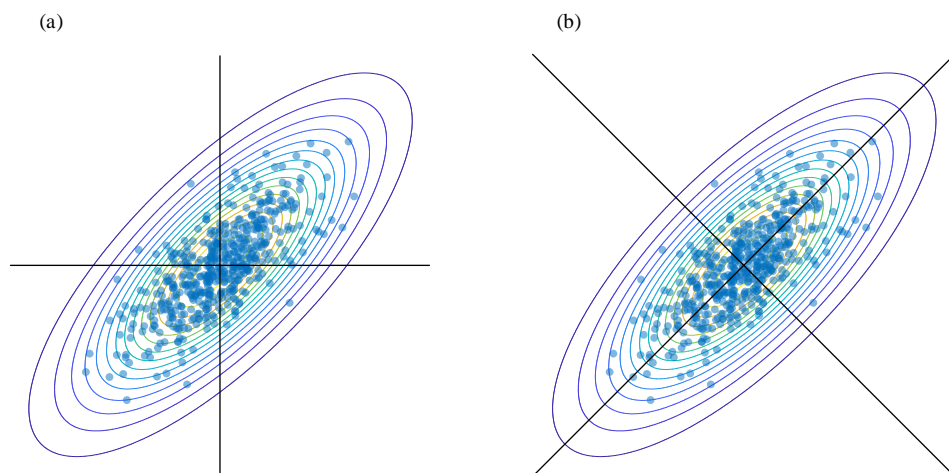


图 18. 两个角度看数据

## 举个例子

图 19 所示为原始二维数据  $\mathbf{X}$ ，可以发现数据的质心位于  $[1, 2]^T$ 。分析数据  $\mathbf{X}$ ，可以发现数据的两个特征上分布分散情况相似，也就是方差大小几乎相同。

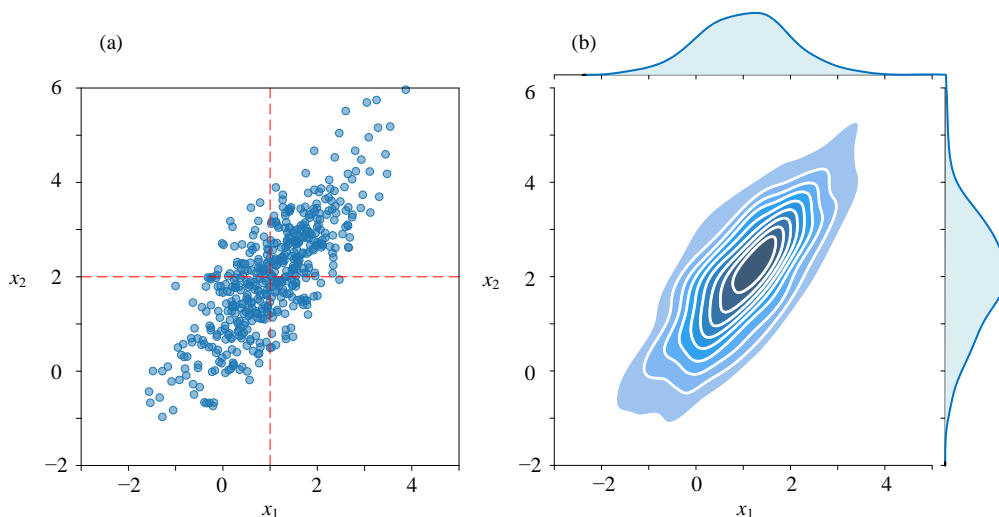


图 19. 原始二维数据  $\mathbf{X}$

利用 `sklearn.decomposition.PCA()` 函数，我们可以通过 `pca.components_` 获得主成分向量。利用 `pca.transform(X)` 可以获得投影后的数据  $\mathbf{Y}$ 。图 20 对比  $\mathbf{Y}$  两列数据分布。图 21 所示为数据  $\mathbf{Y}$  在  $[\mathbf{v}_1, \mathbf{v}_2]$  中散点图。

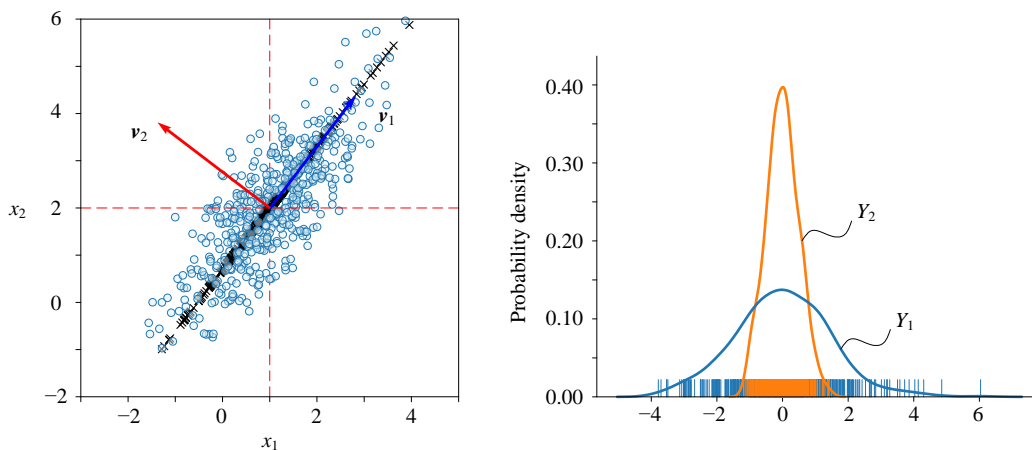
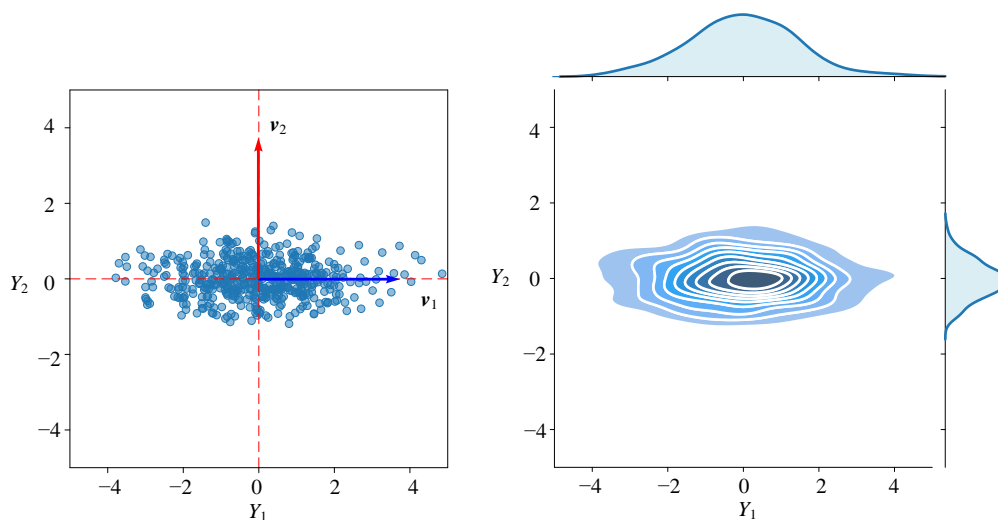


图 20. 主成分数据分布

图 21. 数据  $Y$  在  $[v_1, v_2]$  中散点图

Bk5\_Ch25\_01.py 绘制图 19 ~ 图 21。

## 25.6 奇异值分解

### 四种奇异值分解

丛书在《矩阵力量》一本系统讲解过，奇异值分解有四种形式：

- ◀ **完全型** (full)
- ◀ **经济型** (economy-size, thin)
- ◀ **紧凑型** (compact)
- ◀ **截断型** (truncated)

如图 22 所示，完全型奇异值分解  $S$  矩阵并非方阵，相当于主成分特征值方阵  $S$  下面叠加一块全 0 矩阵  $O$ 。



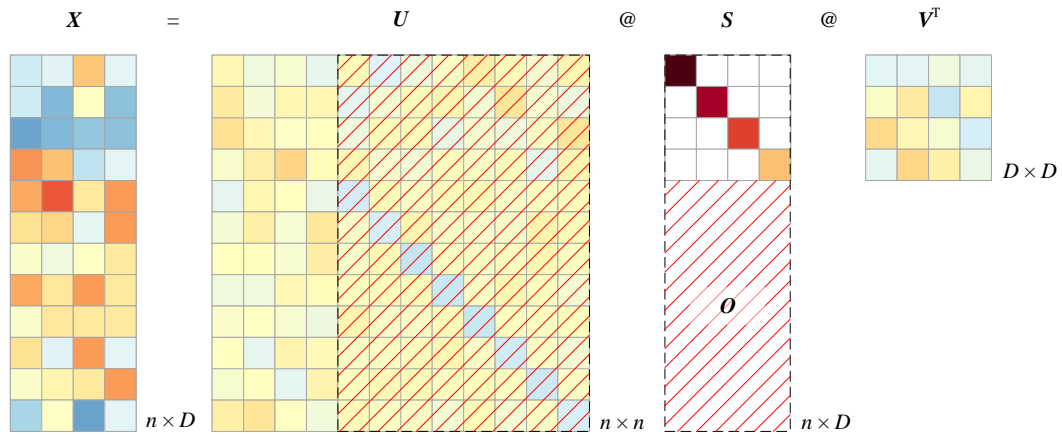


图 22. 完全 (full) 奇异值分解

去掉图 22 中这个全 0 矩阵  $O$ ，便得到经济型奇异值分解。图 23 给出的是经济型奇异值分解， $S$  矩阵为方阵，形状为  $D \times D$ 。

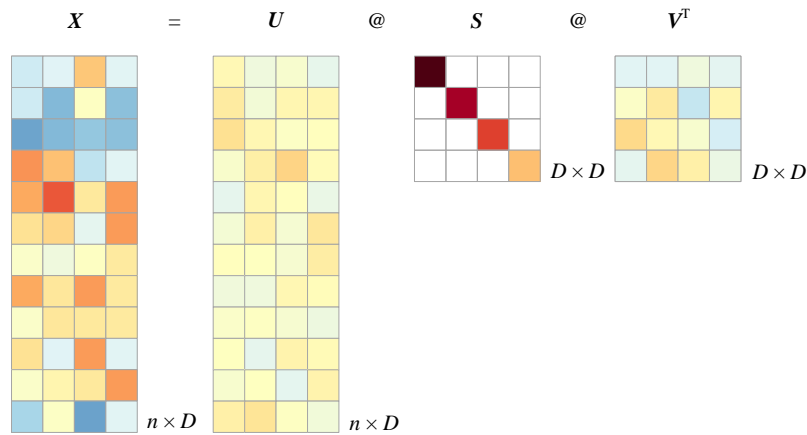


图 23. 经济型奇异值分解

当  $X$  不是满秩时，即  $\text{rank}(X) = r < D$ ，图 23 经济型奇异值分解可以进一步简化为如图 24 所示的紧凑型 SVD 分解。

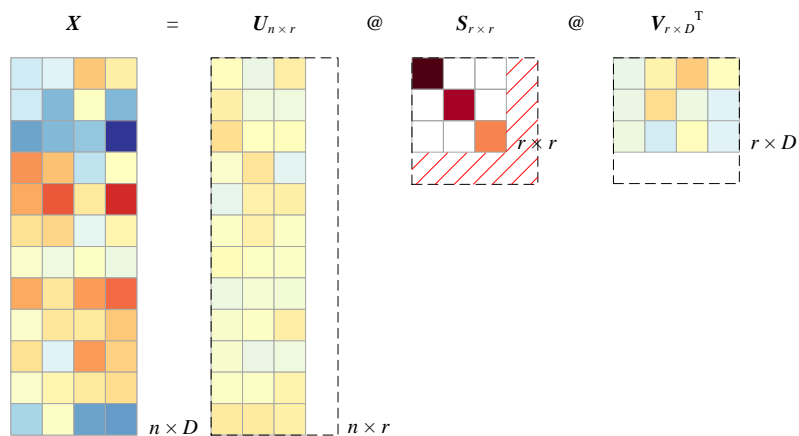


图 24. 紧凑型奇异值分解， $X$  不是满秩

图 25 给出的是截断型奇异值分解， $S_{p \times p}$  仅使用图 23 中  $S$  矩阵  $p$  个主成分特征值，形状为  $p \times p$ 。注意，图 25 中使用的是约等号“ $\approx$ ”；这是因为，约等号右侧矩阵运算仅仅还原  $X$  矩阵部分数据，并非还原全部信息。本章后续将会展开讲解数据还原和误差。

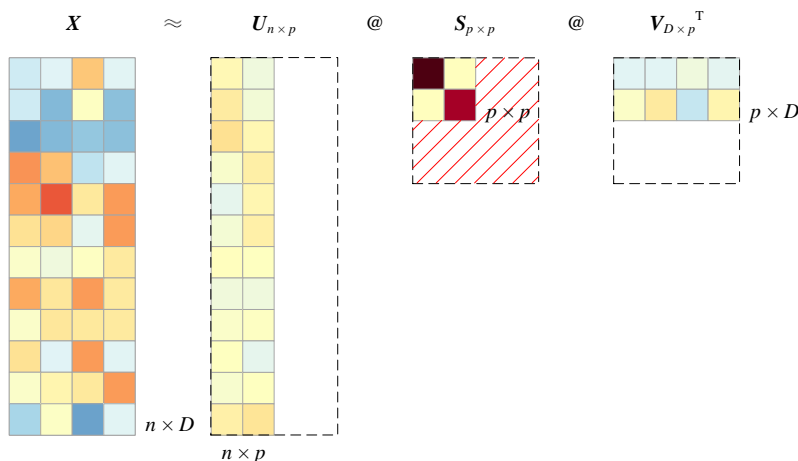


图 25. 截断型奇异值分解

## SVD 完成主成分分析

**奇异值分解** (singular value decomposition, SVD) 也可以用来做主成分分析。首先中心化 (去均值) 数据矩阵。对已经去均值的矩阵  $X_{n \times D}$  进行完全型 SVD 分解，得到：

$$X = USV^T \quad (23)$$

$V$  和  $U$  均为正交矩阵，即满足：

$$\begin{aligned} V^T V &= I \\ U^T U &= I \end{aligned} \quad (24)$$

Python 中常用奇异值分解函数为 `numpy.linalg.svd()`。

由于  $X$  已经中心化，其协方差矩阵可以通过下式计算获得：

$$\Sigma = \frac{X^T X}{n-1} \quad (25)$$

将(23) 代入 (25) 得到：

$$\Sigma = \frac{(USV^T)^T USV^T}{n-1} = \frac{VS^T SV^T}{n-1} \quad (26)$$

对协方差矩阵进行特征值分解：

$$\Sigma = \Lambda V^T \quad (27)$$

联立 (26) 和 (27),

$$\frac{\mathbf{V}\mathbf{S}^T\mathbf{S}\mathbf{V}^T}{n-1} = \mathbf{V}\mathbf{A}\mathbf{V}^T \quad (28)$$

对于经济型 SVD 分解,  $\mathbf{S}$  为对角方阵, (28) 整理得到:

$$\frac{\mathbf{S}^2}{n-1} = \mathbf{A} \quad (29)$$

即

$$\frac{1}{n-1} \begin{bmatrix} s_1^2 & & & \\ & s_2^2 & & \\ & & \ddots & \\ & & & s_D^2 \end{bmatrix} = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_D \end{bmatrix} \quad (30)$$

注意,  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D$ 。

奇异值和特征值存在如下关系:

$$\frac{s_j^2}{n-1} = \lambda_j \quad (31)$$

$s_j$  为第  $j$  个主成分的**奇异值** (singular value),  $\lambda_j$  为协方差矩阵的第  $j$  个特征值。

## 理解 $\mathbf{U}$

$\mathbf{Z}$  反向可以还原  $\mathbf{X}$ :

$$\mathbf{X} = \mathbf{Z}\mathbf{V}^{-1} = \mathbf{Z}\mathbf{V}^T \quad (32)$$

对比 (23) 和  $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ , 可以发现:

$$\mathbf{Z} = \mathbf{U}\mathbf{S} \quad (33)$$

也就是

$$\begin{bmatrix} \mathbf{z}_1 & \mathbf{z}_2 & \cdots & \mathbf{z}_D \end{bmatrix} = \begin{bmatrix} \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_D \end{bmatrix} \begin{bmatrix} s_1 & & & \\ & s_2 & & \\ & & \ddots & \\ & & & s_D \end{bmatrix} = \begin{bmatrix} s_1\mathbf{u}_1 & s_2\mathbf{u}_2 & \cdots & s_D\mathbf{u}_D \end{bmatrix} \quad (34)$$

即:

$$s_1\mathbf{u}_1 = \mathbf{z}_1, \quad s_2\mathbf{u}_2 = \mathbf{z}_2, \quad \dots \quad (35)$$

对  $\mathbf{z}_1$  求方差:

$$\text{var}(\mathbf{z}_1) = \frac{\mathbf{z}_1^T \mathbf{z}_1}{n-1} = \frac{(s_1\mathbf{u}_1)^T (s_1\mathbf{u}_1)}{n-1} = \frac{s_1^2 \|\mathbf{u}_1\|^2}{n-1} = \frac{s_1^2}{n-1} = \lambda_1 \quad (36)$$

可以发现矩阵  $U$  每一列数据相当于  $Z$  的标准化：

$$U = [u_1 \quad u_2 \quad \cdots \quad u_D] = \begin{bmatrix} \frac{z_1}{s_1} & \frac{z_2}{s_2} & \cdots & \frac{z_D}{s_D} \end{bmatrix} \quad (37)$$

也就是：

$$U = [u_1 \quad u_2 \quad \cdots \quad u_D] = ZS^{-1} \quad (38)$$

至此，我们理解了 SVD 分解中矩阵  $U$  的内涵。

## 张量积

用张量积来展开 SVD 分解：

$$\begin{aligned} X &= USV^T \\ &= [u_1 \quad u_2 \quad \cdots \quad u_D] \begin{bmatrix} s_1 & & & \\ & s_2 & & \\ & & \ddots & \\ & & & s_D \end{bmatrix} \begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_D^T \end{bmatrix} \\ &= s_1 u_1 v_1^T + s_2 u_2 v_2^T + \cdots + s_D u_D v_D^T \\ &= s_1 u_1 \otimes v_1 + s_2 u_2 \otimes v_2 + \cdots + s_D u_D \otimes v_D \end{aligned} \quad (39)$$

图 26 所示为 (39) 还原原始数据的过程。

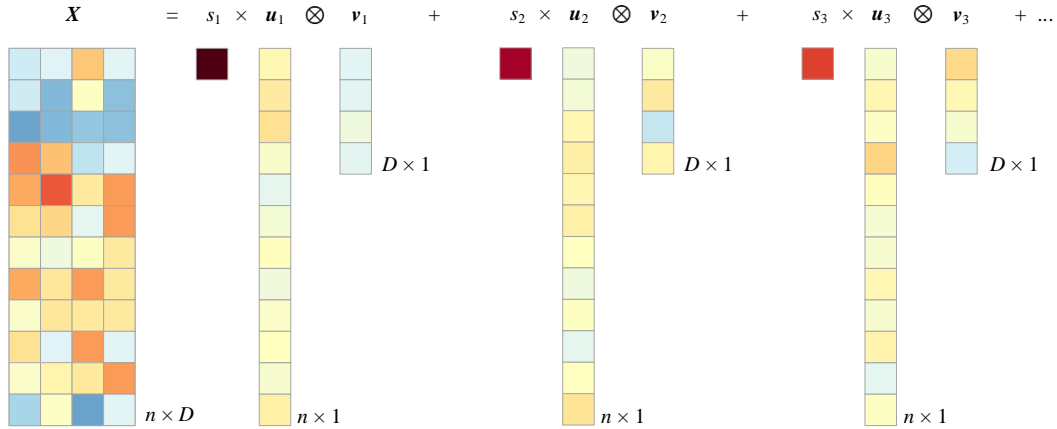


图 26. 张量积  $s_1 u_1 \otimes v_1$ 、 $s_2 u_2 \otimes v_2$  等之和还原数据  $X$

## 25.7 优化问题

本节最后要从两个角度构造主成分分析优化问题角度，并剖析主成分分析理论内核。如图 27 所示， $\mathbf{X}$  为中心化数据，数据中心为零向量；数据  $\mathbf{X}$  协方差矩阵如下：

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{n-1} \quad (40)$$

图 27 中， $\mathbf{v}$  为某个主成分向量。数据  $\mathbf{X}$  在  $\mathbf{v}$  上投影结果为  $\mathbf{z}$ ：

$$\mathbf{z} = \mathbf{X}\mathbf{v} \quad (41)$$

主成分分析中，选取主成分  $\mathbf{v}$  向量核心是在方向上数据投影值  $\mathbf{z}$  方差最大化；这便是构造优化问题第一个角度。

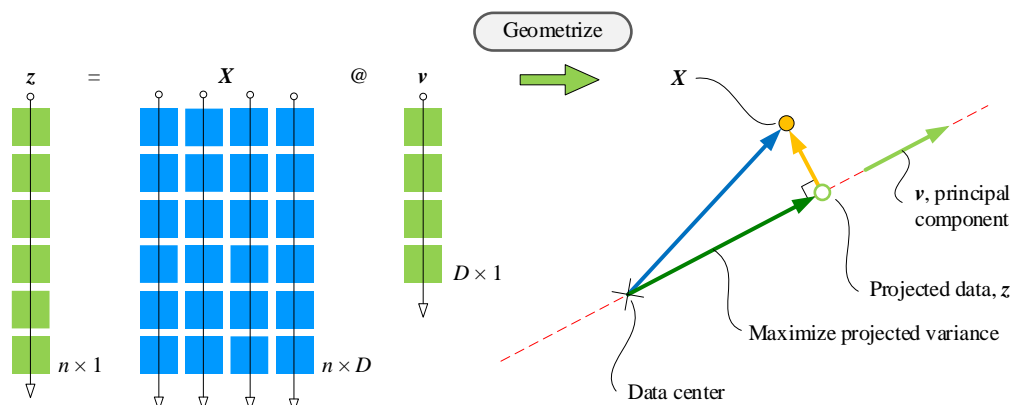


图 27. 主成分分析优化问题

由于  $\mathbf{X}$  为中心化数据，因此  $\mathbf{z}$  的均值也为 0；因此， $\mathbf{z}$  方差为：

$$\text{var}(\mathbf{z}) = \frac{\mathbf{z}^T \mathbf{z}}{n-1} = \mathbf{v}^T \overset{\text{Covariance matrix}}{\frac{\mathbf{X}^T \mathbf{X}}{n-1}} \mathbf{v} \quad (42)$$

发现上式隐藏着数据  $\mathbf{X}$  协方差矩阵，因此  $\text{var}(\mathbf{z})$  为：

$$\text{var}(\mathbf{z}) = \mathbf{v}^T \Sigma \mathbf{v} \quad (43)$$

$\mathbf{v}$  为单位列向量，即下式成立：

$$\mathbf{v}^T \mathbf{v} = 1 \quad (44)$$

有以上分析，构造主成分分析优化问题，优化目标为数据在  $\mathbf{v}$  方向上数据投影值方差最大化：

$$\begin{aligned} \arg \max_{\mathbf{v}} \quad & \mathbf{v}^T \Sigma \mathbf{v} \\ \text{subject to:} \quad & \mathbf{v}^T \mathbf{v} - 1 = 0 \end{aligned} \quad (45)$$

上式最大化优化问题等价于如下最小化优化问题：

$$\begin{aligned} \arg \min_{\mathbf{v}} \quad & -\mathbf{v}^T \Sigma \mathbf{v} \\ \text{subject to: } & \mathbf{v}^T \mathbf{v} - 1 = 0 \end{aligned} \quad (46)$$

构造拉格朗日函数  $L(\mathbf{v}, \lambda)$ ：

$$L(\mathbf{v}, \lambda) = -\mathbf{v}^T \Sigma \mathbf{v} + \lambda (\mathbf{v}^T \mathbf{v} - 1) \quad (47)$$

$\lambda$  为拉格朗日乘子。 $L(\mathbf{x}, \lambda)$  对  $\mathbf{v}$  求偏导，最优解必要条件如下：

$$\nabla_{\mathbf{v}} L(\mathbf{v}, \lambda) = \frac{\partial L(\mathbf{v}, \lambda)}{\partial \mathbf{v}} = (-2\Sigma \mathbf{v} + 2\lambda \mathbf{v})^T = \mathbf{0} \quad (48)$$

有关拉格朗日乘子法，请大家回顾《矩阵力量》第 18 章。

整理 (48) 得到：

$$\Sigma \mathbf{v} = \lambda \mathbf{v} \quad (49)$$

由此， $\mathbf{v}$  为数据  $\mathbf{X}$  方差-协方差矩阵  $\Sigma$  特征向量。 $\text{var}(\mathbf{z})$  整理为：

$$\text{var}(\mathbf{z}) = \mathbf{v}^T \Sigma \mathbf{v} = \mathbf{v}^T \lambda \mathbf{v} = \lambda \mathbf{v}^T \mathbf{v} = \lambda \quad (50)$$

即说， $\text{var}(\mathbf{z})$  最大值对应  $\Sigma$  最大特征值。 $\Sigma$  特征值分解：

$$\Sigma = \mathbf{V} \mathbf{A} \mathbf{V}^T \quad (51)$$

其中，

$$\mathbf{V} = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \cdots \quad \mathbf{v}_d] \\ \mathbf{A} = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_d \end{bmatrix}, \quad \lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d \quad (52)$$

这一节从优化角度解释了为什么特征值分解能够完成主成分分析。

## 25.8 数据还原和误差

### 还原

前文介绍过， $\mathbf{Z}$  反向可以通过  $\mathbf{X} = \mathbf{Z} \mathbf{V}^T$  还原  $\mathbf{X}$ 。

图 28 所示为还原得到  $\mathbf{X}$  过程。

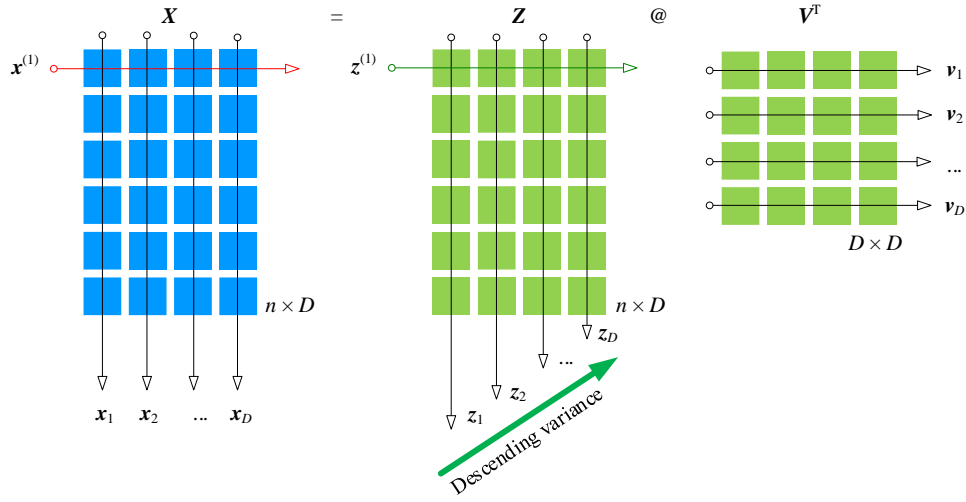


图 28. 反向还原数据  $X = ZV^T$

图 29 所示热图，为新特征数据矩阵  $Z$  还原转化为原始数据矩阵  $X$ 。

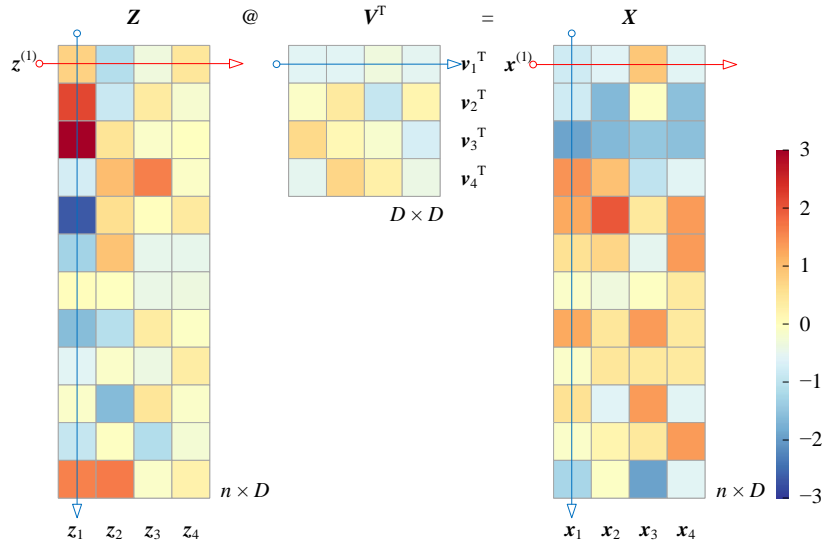


图 29. 新特征数据矩阵  $Z$  还原转化为原始数据矩阵  $X$

$X = ZV^T$  展开得到下式：

$$X = \begin{bmatrix} z_1 & z_2 & z_3 & z_4 \end{bmatrix} \begin{bmatrix} v_1^T \\ v_2^T \\ v_3^T \\ v_4^T \end{bmatrix} = \begin{bmatrix} z_1 v_1^T + z_2 v_2^T + z_3 v_3^T + z_4 v_4^T \\ \hat{x}_1 & \hat{x}_2 & \hat{x}_3 & \hat{x}_4 \end{bmatrix} \quad (53)$$

(53) 所示运算过程如图 30 所示。

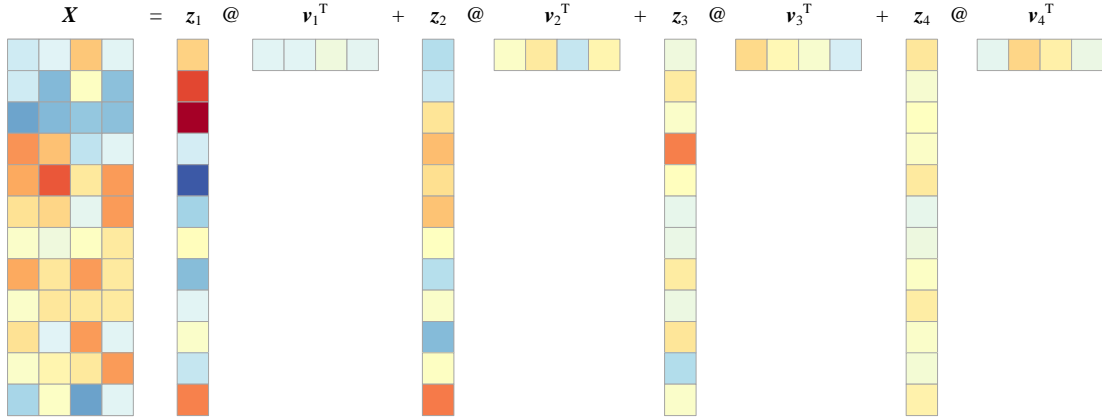


图 30. 还原原始数据运算

图 31 所示为  $z_1$  还原  $X$  部分数据，对应运算如下：

$$X_1 = z_1 v_1^T \quad (54)$$

展开上式得到：

$$\begin{aligned} X_1 &= z_1 v_1^T \\ &= z_1 \begin{bmatrix} v_{1,1} & v_{2,1} & \cdots & v_{D,1} \end{bmatrix} \\ &= \begin{bmatrix} v_{1,1} z_1 & v_{2,1} z_1 & \cdots & v_{D,1} z_1 \end{bmatrix} \end{aligned} \quad (55)$$

观察图 31 热图可以发现一些有意思的特点。还原得到的数据每一列热图模式高度相似；(55) 解释了这一点， $X_1$  的每一列均是标量乘以向量  $z_1$  的结果。显然， $X_1$  的秩为 1，即  $\text{rank}(X_1) = 1$ 。

图 32、图 33 和图 34 分别展示  $z_2$ 、 $z_3$  和  $z_4$  还原  $X$  部分数据。

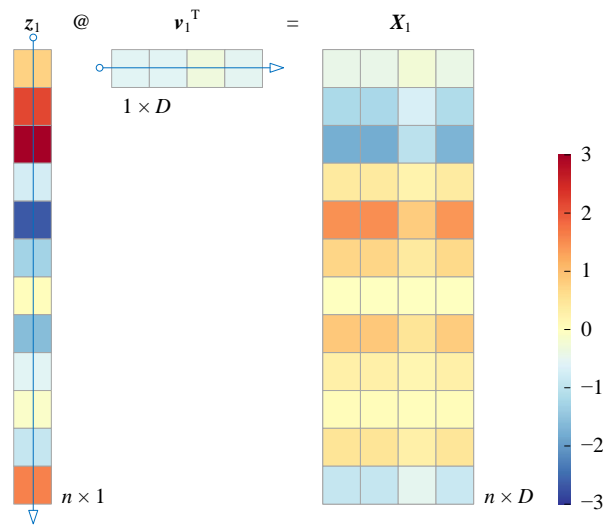


图 31.  $z_1$  还原  $X$  部分数据



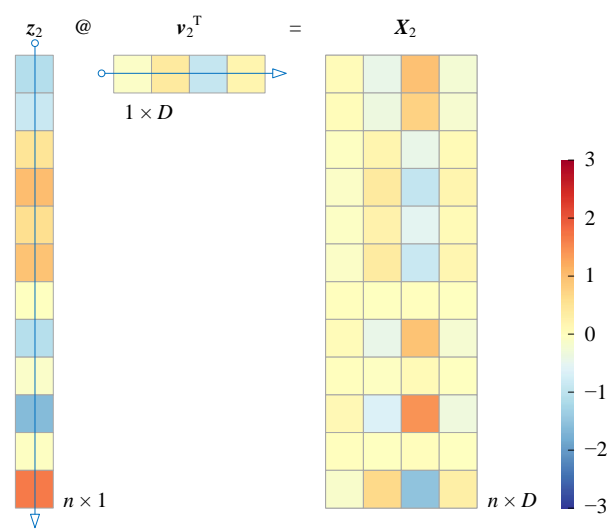


图 32.  $z_2$  还原  $X$  部分数据

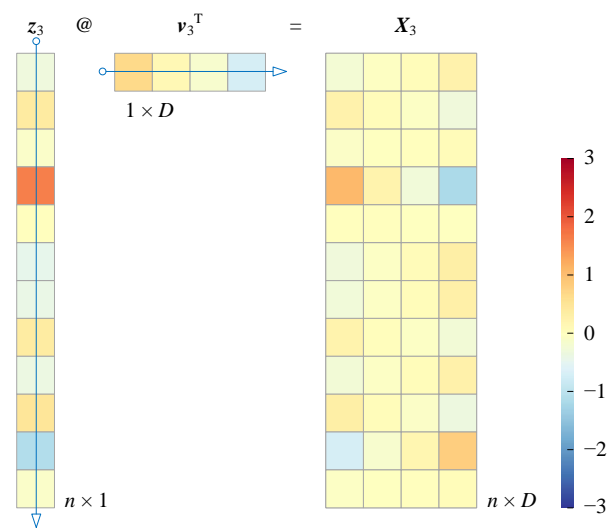


图 33.  $z_3$  还原  $X$  部分数据

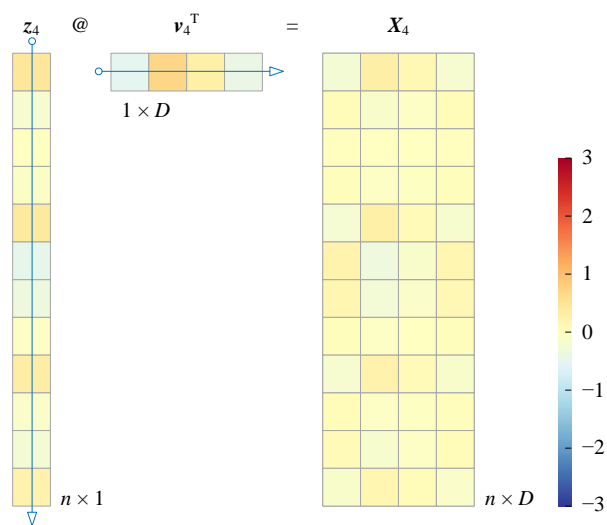
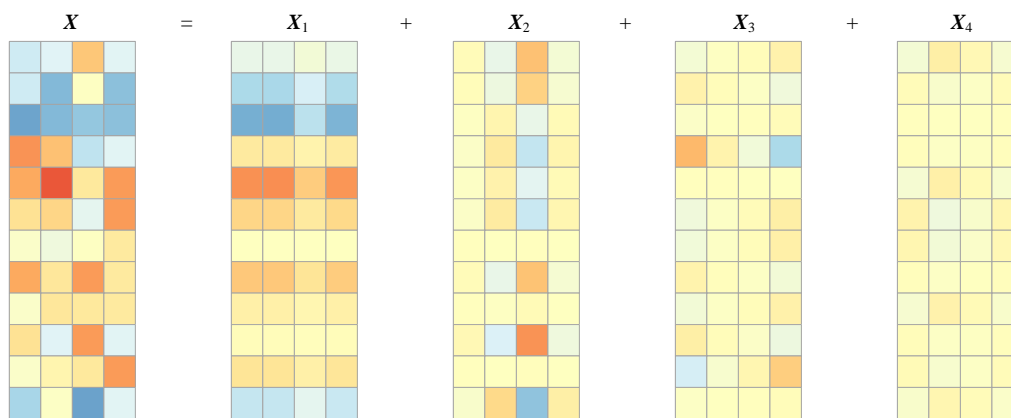
图 34.  $z_4$  还原  $X$  部分数据

图 35 所示为原始数据矩阵  $X$  热图于四层热图叠加结果。观察图 35，发现随着主成分次数降低，每个主成分各自对数据  $X$  还原力度不断降低，看到还原热图颜色越来越浅；但是，把这些主成分各自还原生成热图不断叠加，获得热图就不断逼近原始热图。

图 35. 原始数据矩阵  $X$  热图于四层热图叠加结果

## 张量积

另外，(53) 可以用张量积来表达：

$$X = \underbrace{z_1 \otimes v_1}_{\hat{X}_1} + \underbrace{z_2 \otimes v_2}_{\hat{X}_2} + \underbrace{z_3 \otimes v_3}_{\hat{X}_3} + \underbrace{z_4 \otimes v_4}_{\hat{X}_4} \quad (56)$$

图 36 所示为通过主成分  $v_1, v_2, v_3, v_4$  和其自身转置乘积计算张量积。

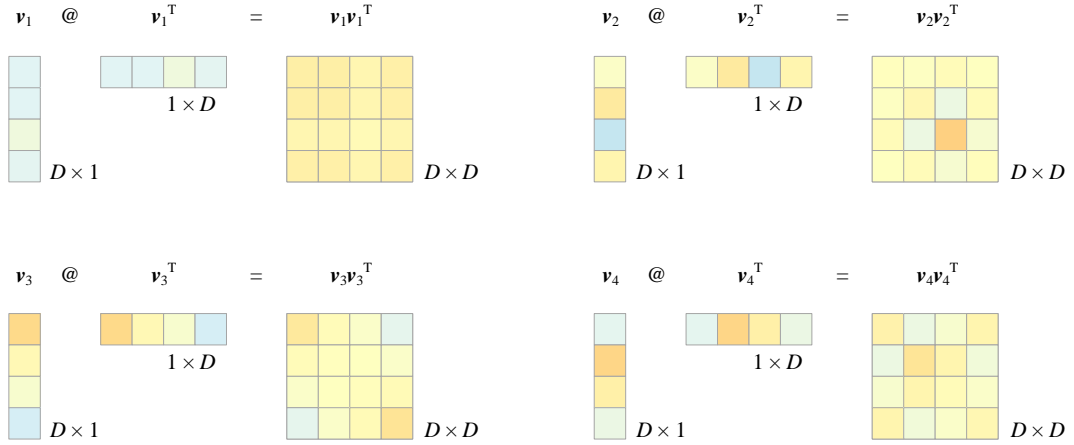


图 36. 列向量乘自身转置获得四个张量积

图 37 所示为张量积运算，和图 36 结果完全一致。

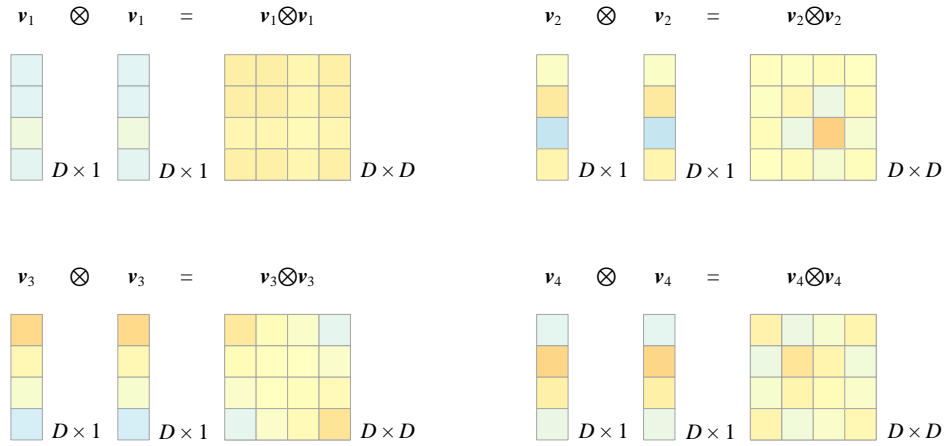


图 37. 内积计算获得四个张量积

利用 (14), (53) 可以整理为：

$$X = Xv_1v_1^T + Xv_2v_2^T + \dots + Xv_Dv_D^T = \sum_{j=1}^D Xv_jv_j^T = X \left( \sum_{j=1}^D v_jv_j^T \right) \quad (57)$$

(57) 可以用张量积表达：

$$X = X(v_1 \otimes v_1) + X(v_2 \otimes v_2) + \dots + X(v_D \otimes v_D) = \sum_{j=1}^D Xv_j \otimes v_j = X \left( \sum_{j=1}^D v_j \otimes v_j \right) \quad (58)$$

容易推导得到，(58) 中张量积相加得到单位矩阵。《矩阵力量》第 10 章给这种投影一个特别的名字——二次投影，建议大家回顾。

$$\mathbf{v}_1 \otimes \mathbf{v}_1 + \mathbf{v}_2 \otimes \mathbf{v}_2 + \dots + \mathbf{v}_D \otimes \mathbf{v}_D = \left( \sum_{j=1}^D \mathbf{v}_j \otimes \mathbf{v}_j \right) = \mathbf{I} \quad (59)$$

上式如图 38 热图所示。

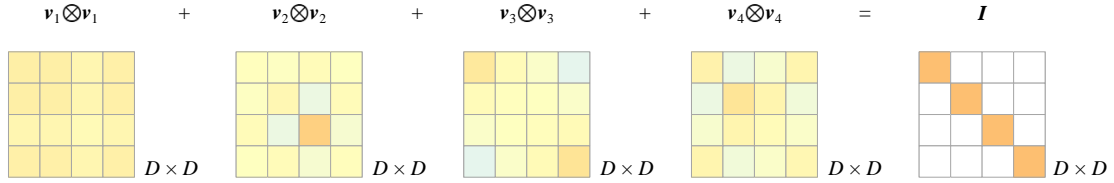


图 38. 张量积相加得到单位矩阵

联立 (16) 和 (54)，利用张量积  $\mathbf{v}_1 \otimes \mathbf{v}_1$  还原部分原始数据：

$$\mathbf{X}_1 = \mathbf{z}_1 \mathbf{v}_1^T = \mathbf{X} \mathbf{v}_1 \mathbf{v}_1^T = \mathbf{X} \underbrace{(\mathbf{v}_1 \otimes \mathbf{v}_1)}_{\text{Tensor product}} \quad (60)$$

类似，张量积  $\mathbf{v}_2 \otimes \mathbf{v}_2$  也可以还原部分原始数据：

$$\mathbf{X}_2 = \mathbf{z}_2 \mathbf{v}_2^T = \mathbf{X} \mathbf{v}_2 \mathbf{v}_2^T = \mathbf{X} \underbrace{(\mathbf{v}_2 \otimes \mathbf{v}_2)}_{\text{Tensor product}} \quad (61)$$

图 39 所示为张量积  $\mathbf{v}_1 \otimes \mathbf{v}_1$  和  $\mathbf{v}_2 \otimes \mathbf{v}_2$  还原部分数据  $\mathbf{X}$ ；图 40 所示为张量积  $\mathbf{v}_3 \otimes \mathbf{v}_3$  和  $\mathbf{v}_4 \otimes \mathbf{v}_4$  还原部分数据  $\mathbf{X}$ 。

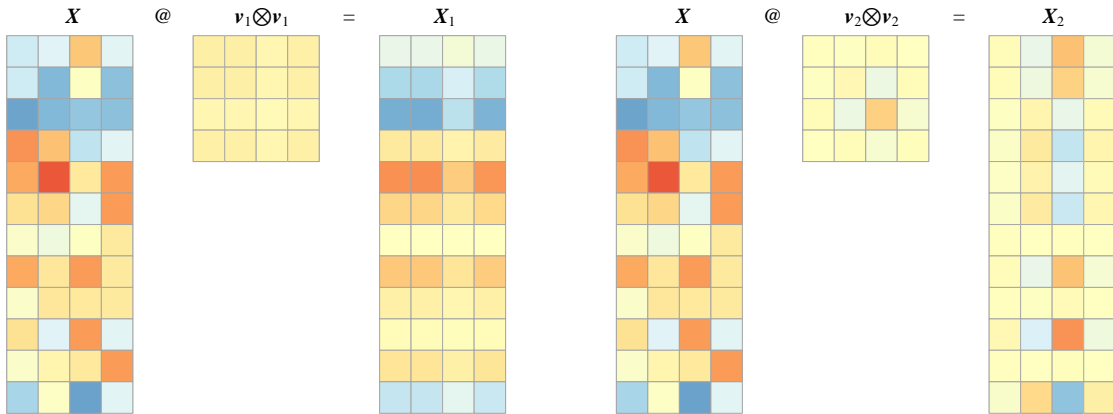
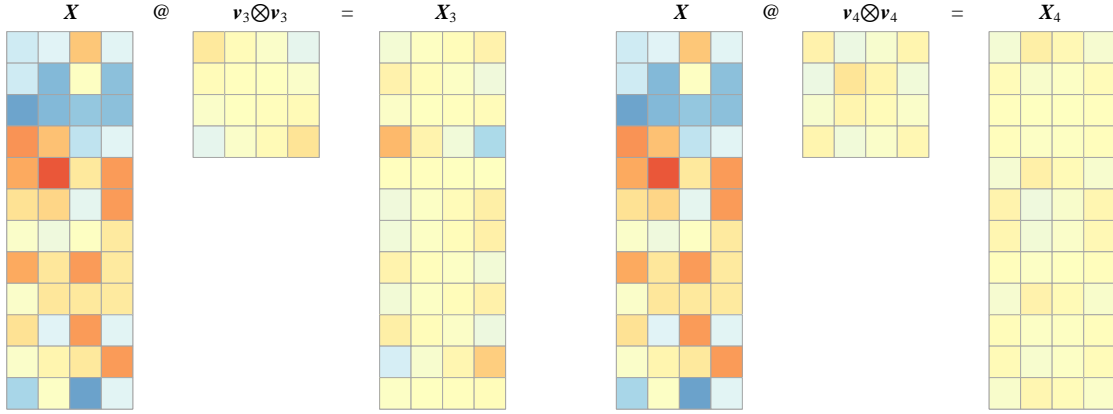


图 39. 张量积  $\mathbf{X}(\mathbf{v}_1 \otimes \mathbf{v}_1)$  和  $\mathbf{X}(\mathbf{v}_2 \otimes \mathbf{v}_2)$  还原部分数据  $\mathbf{X}$

图 40. 张量积  $X(v_3 \otimes v_3)$  和  $X(v_4 \otimes v_4)$  还原部分数据  $X$ 

## 误差

图 41 所示为两个主成分  $v_1$  和  $v_2$  还原获得原始数据热图，具体计算如下：

$$\hat{X} = [z_1 \ z_2] [v_1 \ v_2]^T \quad (62)$$

相当于

$$\begin{aligned} \hat{X} &= X_1 + X_2 = z_1 v_1^T + z_2 v_2^T \\ &= X(v_1 v_1^T + v_2 v_2^T) = X(v_1 \otimes v_1 + v_2 \otimes v_2) \end{aligned} \quad (63)$$

图 42 所示为通过叠加图 31 和图 32 两个热图还原原始数据矩阵。

从张量积角度来看图 42,

$$X \approx s_1 u_1 \otimes v_1 + s_2 u_2 \otimes v_2^T \quad (64)$$

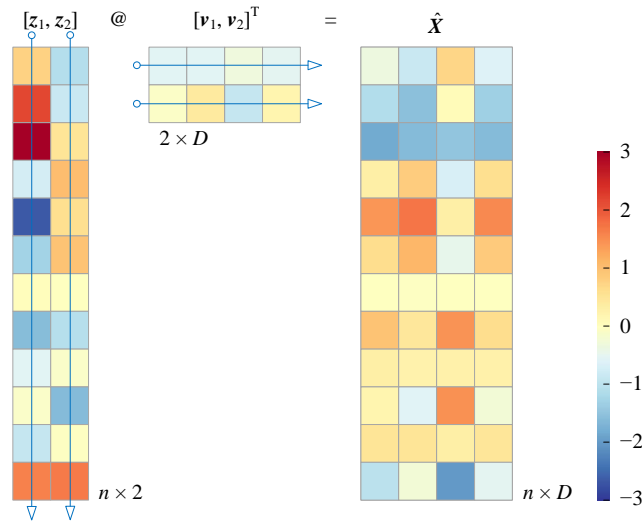


图 41. 前两个主成分  $z_1$  和  $z_2$  还原  $X$  数据

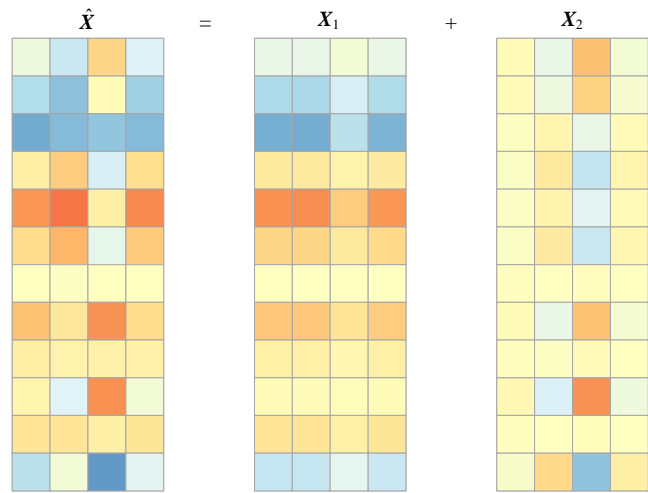


图 42. 两个热图叠加还原原始数据

残差数据矩阵  $E$ ，即原始热图和还原热图色差，利用下式计算获得：

$$E = X - \hat{X}$$

(65)

图 43 比较原始数据  $X$ 、拟合数据  $\hat{X}$  和残差数据矩阵  $E$  热图，发现原始数据  $X$  和拟合数据  $\hat{X}$  已经相差无几。从图片还原角度，如图 43 所示，PCA 降维用更少维度、更少数据获得几乎一样画质图片。

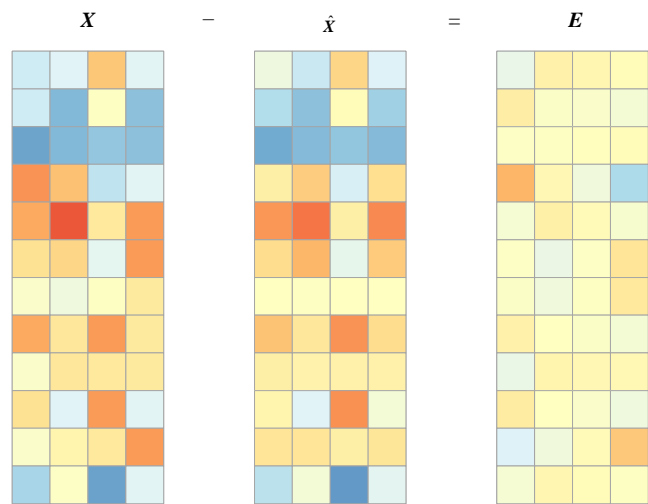


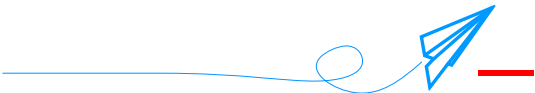
图 43. 原始数据、拟合数据和残差数据热图

## 六条技术路径

相信大家对表 1 并不陌生，大家都在《矩阵力量》第 25 章中见过这六条 PCA 技术路线。本章介绍的实际上是：a) 特征值分解协方差矩阵；b) 奇异值分解中心化数据矩阵。《数据有道》一册将比较表 1 这六种方法的异同。

表 1. 六条 PCA 技术路线，来自《矩阵分解》第 25 章

对象	方法	结果
原始数据矩阵 $X$	奇异值分解	$X = U_X S_X V_X^T$
格拉姆矩阵 $G = X^T X$ 本章中用“修正”的格拉姆矩阵 $G = \frac{X^T X}{n-1}$	特征值分解	$G = V_X \Lambda_X V_X^T$
中心化数据矩阵 $X_c = X - E(X)$	奇异值分解	$X_c = U_c S_c V_c^T$
协方差矩阵 $\Sigma = \frac{(X - E(X))^T (X - E(X))}{n-1}$	特征值分解	$\Sigma = V_c \Lambda_c V_c^T$
标准化数据 (z 分数) $Z_X = (X - E(X)) D^{-1}$ $D = \text{diag}(\text{diag}(\Sigma))^{\frac{1}{2}}$	奇异值分解	$Z_X = U_Z S_Z V_Z^T$
相关性系数矩阵 $P = D^{-1} \Sigma D^{-1}$ $D = \text{diag}(\text{diag}(\Sigma))^{\frac{1}{2}}$	特征值分解	$P = V_Z \Lambda_Z V_Z^T$



人类的思维方式天然具备概率统计属性。概率统计的背后的思想更贴近“生活常识”。人们在谈论可能性的时候，大脑就不自觉进入“概率统计”模式。看着天上云层很厚，可能两小时就会下雨。昨晚淋了雨，估计今天要感冒。估计这次考试通过率 80% 以上。咱们剪刀石头布，三局两胜决胜负。

可惜的是，当数学家将这些生活常识“翻译成”数学语言之后，它们就变成了冷冰冰“火星文”。更遗憾的是，很多概率统计图书读起来更像是数学公式手册，满篇的公式、定理、推导。

概率统计与其说是工具，不如说是方法论、世界观。大家常说的“一命，二运，三风水，四读书”，体现的也是概率统计的思维。

命运不可问，命中没有莫强求。“小概率事件”能发生，得之我幸，不得我命。

风水轮流转，玄而又玄。

正所谓知识改变命运，只有读书成才对应“大概率事件”。大家捧起这本书的时候，就依靠统计思维做出了“最优化”选择。

《统计至简》是“数学”板块的三本中的最后一本。大家读到这里，也就完成了整个“数学”板块的修炼。希望大家日后在看到任何公式的时候，闭上眼睛，能够在脑海中“看见”矩阵和各种几何图形。

下面，我们将踏上《数据有道》、《机器学习》的“实践”之旅！期待和大家共同学习、成长！