

香港中文大學
數學系
主成分分析練習

1 引言(Introduction)

主成分分析(Principal Component Analysis, PCA) 最早由卡爾·皮爾遜(Karl Pearson) 於1901 年正式提出，隨後哈羅德·霍特林(Harold Hotelling) 於1933 年獨立地重新發展並擴展了該方法。皮爾遜指出，他最初的代數框架可輕易應用於一般的數值數據集。儘管他承認對於具有四個或以上不同變數的數據集，人手計算會變得繁重，但他證實了對於小規模問題，完全由人手計算在數學上是可行的。人手計算PCA 在實際操作中僅適用於最多四個變數的數據集；然而，PCA 正正是在應用於具有五個或以上高維特徵的數據集時，才能發揮其最大的分析與計算優勢。

嚴格落實PCA 需要具備描述性統計學及線性矩陣代數的堅實基礎，包括向量點積(vector dot products)、內積(inner products)、在直線及高維子空間上的正交投影(orthogonal projection)、特徵分解(eigendecomposition) 以及協方差矩陣代數(covariance matrix algebra)。

PCA 的核心概念目標具有兩個數學上等價的表述方式：

1. 最大化零中心化數據(mean-centered data) 投影在單位基底向量上的總投影變異數(variance)
2. 最小化原始數據點與其在單位基底向量上投影點之間的總正交距離平方(squared orthogonal distance)

PCA 的實際動機源於未經處理的高維數據本身固有的複雜性：冗餘及相關的特徵會增加計算運行時間並引入統計上的不穩定性。在數學上，主成分是透過對原始特徵集進行正交線性轉換來構建的，其中轉換矩陣優化了一個「變異數最大化」的代數目標準則。PCA 被歸類為一種無監督學習(unsupervised learning) 算法：在轉換過程中不需要標記輸出/目標響應變數，其核心優化目標是最大化保留在投影子空間中的總可解釋變異數(explained variance)。

2 PCA 演算法流程(PCA Algorithm Pipeline)

執行PCA 的完整順序工作流程定義為以下五個明確的有序階段：

1. 特徵標準化/ 平均值中心化（數據縮放）
2. 計算中心化數據的樣本協方差矩陣
3. 協方差矩陣的特徵值與特徵向量分解
4. 按特徵值大小降序排列特徵向量，構建主成分基底
5. 將原始中心化數據集投影到截斷的特徵向量基底上，生成降維後的輸出數據集

3 PCA 的基礎線性代數：點積、內積、範數、直線上投影

推導PCA 所需的所有向量運算均在下方提供了逐步的代數推導與完整定義。

3.1 點積 / 標準歐幾里得內積(Dot Product)

點積是實歐幾里得向量空間 \mathbb{R}^n 中的典範內積。它接受兩個等長的向量並返回一個純量實數。

3.1.1 定義

設 $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ 為兩個 n 維實向量，其元素為 x_i, y_i ($i = 1, \dots, n$)。點積寫為 \mathbf{x} 的轉置矩陣與 \mathbf{y} 的矩陣乘法：

$$\mathbf{x}^\top \mathbf{y} = \sum_{i=1}^N x_i y_i \quad (3.1)$$

3.1.2 基於點積的正交條件

當且僅當兩個向量的點積為零時，這兩個向量被稱為正交（幾何上垂直）：

$$\mathbf{x}^\top \mathbf{y} = 0 \quad (3.2)$$

3.2 向量範數 / 向量長度(Vector Norm)

向量的歐幾里得範數（長度）量化了其在歐幾里得空間中的大小，透過向量與自身的點積來定義。

3.2.1 定義

對於任何向量 $\mathbf{x} \in \mathbb{R}^n$ ：

$$\|\mathbf{x}\| = \sqrt{\mathbf{x}^\top \mathbf{x}} = \sqrt{\sum_{i=1}^N x_i^2} \quad (3.3)$$

3.2.2 L^2 範數

將範數平方即可得出各元素的平方和：

$$\|\mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{x} = \sum_{i=1}^N x_i^2$$

3.3 兩個向量之間的歐幾里得距離

距離量度的是 \mathbf{x} 與 \mathbf{y} 之間的差向量的大小。

$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{(\mathbf{x} - \mathbf{y})^\top (\mathbf{x} - \mathbf{y})} \quad (3.4)$$

3.4 兩個向量之間的夾角

向量 \mathbf{x}, \mathbf{y} 之間幾何夾角 α 的餘弦值與其點積及範數有關：

$$\cos \alpha = \frac{\mathbf{x}^\top \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \quad (3.5)$$

3.5 一般內積空間（點積的推廣）

點積是針對任意實向量空間 V 定義的一種更廣泛的抽象運算（稱為內積）的特例。內積接受兩個向量並返回一個純量，將幾何上的點積行為推廣至非歐幾里得空間。

3.5.1 正式的公理定義

設 V 為在 \mathbb{R} 上的向量空間。內積 $(\cdot, \cdot) : V \times V \rightarrow \mathbb{R}$ 必須滿足四個公理：

1. 對稱性(Symmetry): $(\mathbf{u}, \mathbf{v}) = (\mathbf{v}, \mathbf{u}), \quad \forall \mathbf{u}, \mathbf{v} \in V$
2. 正定性(Positive Definiteness): 對於所有 $\mathbf{u} \neq \mathbf{0}$, $(\mathbf{u}, \mathbf{u}) > 0$; 且 $(\mathbf{0}, \mathbf{0}) = 0$
3. 第一變元線性(Linearity in the first argument): $(\alpha\mathbf{u} + \mathbf{w}, \mathbf{v}) = \alpha(\mathbf{u}, \mathbf{v}) + (\mathbf{w}, \mathbf{v}), \quad \forall \alpha \in \mathbb{R}, \mathbf{u}, \mathbf{v}, \mathbf{w} \in V$
4. 第二變元線性(Linearity in the second argument): $(\mathbf{u}, \alpha\mathbf{v} + \mathbf{w}) = \alpha(\mathbf{u}, \mathbf{v}) + (\mathbf{u}, \mathbf{w}), \quad \forall \alpha \in \mathbb{R}, \mathbf{u}, \mathbf{v}, \mathbf{w} \in V$

3.5.2 基於內積的一般範數

向量大小的概念可以推廣至任何內積空間：

$$\|\mathbf{u}\| = \sqrt{(\mathbf{u}, \mathbf{u})} \quad (3.6)$$

3.5.3 基於內積的一般距離

$$\text{dist}(\mathbf{u}, \mathbf{v}) = \|\mathbf{u} - \mathbf{v}\| = \sqrt{(\mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v})} \quad (3.7)$$

3.5.4 一般正交條件

相對於所選的內積，向量 $\mathbf{u}, \mathbf{v} \in V$ 為正交，當且僅當：

$$(\mathbf{u}, \mathbf{v}) = 0 \quad (3.8)$$

正交性依賴於所使用的具體內積：在某一種內積定義下正交的兩個向量，並不保證在另一種內積定義下依然正交。

3.6 在單一向量直線上的正交投影

PCA 投影的幾何基礎：第一個主成分是唯一的單位向量，當所有數據點投影到該向量方向時，會最大化數據集的總變異數。投影是一種將一個向量映射到由第二個基底向量所展開的直線上的線性算子。

3.6.1 正交投影定義

設 \mathbf{a} 為任意輸入向量，而 \mathbf{b} 為定義目標直線的固定基底向量。 \mathbf{a} 在 \mathbf{b} 上的正交投影是一個恰好位於由 \mathbf{b} 所展開的直線上的向量 $c\mathbf{b}$ ，其中殘差差分向量 $\mathbf{a} - c\mathbf{b}$ 與 \mathbf{b} 正交。殘差 $\mathbf{a} - c\mathbf{b}$ 被稱為正交投影誤差項。這種關係如圖1 所示。

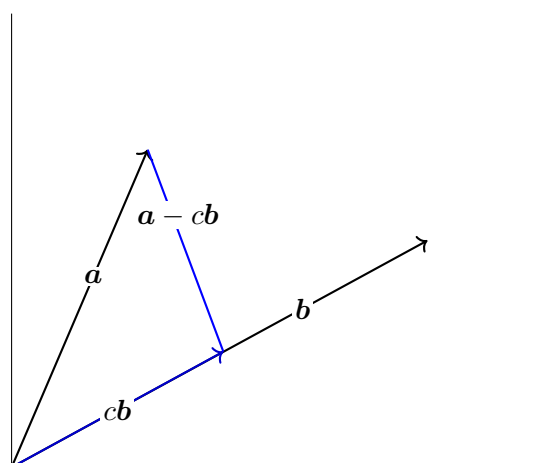


Figure 1: 向量 \mathbf{a} 在由 \mathbf{b} 展開的直線上的正交投影。殘差 $\mathbf{a} - \mathbf{cb} \perp \mathbf{b}$ 。

3.6.2 投影純量 c 的推導

根據正交條件，殘差向量與基底向量的內積為零：

$$\begin{aligned} (\mathbf{a} - \mathbf{cb}, \mathbf{b}) &= 0 \\ (\mathbf{a}, \mathbf{b}) - (\mathbf{cb}, \mathbf{b}) &= 0 \\ (\mathbf{a}, \mathbf{b}) - c(\mathbf{b}, \mathbf{b}) &= 0 \end{aligned}$$

重新排列以分離出純量系數 c ：

$$c(\mathbf{b}, \mathbf{b}) = (\mathbf{a}, \mathbf{b}) \implies c = \frac{(\mathbf{a}, \mathbf{b})}{(\mathbf{b}, \mathbf{b})}$$

對於歐幾里得點積：

$$(\mathbf{b}, \mathbf{b}) = \mathbf{b}^\top \mathbf{b} = \|\mathbf{b}\|^2 \implies c = \frac{\mathbf{b}^\top \mathbf{a}}{\|\mathbf{b}\|^2}$$

3.6.3 單位基底向量的簡化 $\|\mathbf{b}\| = 1$

$$c = \mathbf{b}^\top \mathbf{a}$$

3.6.4 單一向量基底的投影矩陣

投影映射為線性的，其矩陣為 $P_{\mathbf{b}} = \frac{\mathbf{b}\mathbf{b}^\top}{\|\mathbf{b}\|^2}$ ，因此 $\text{proj}_{\mathbf{b}}(\mathbf{a}) = P_{\mathbf{b}}\mathbf{a}$ 。

4 高維子空間上的正交投影

我們將正交投影從一維直線推廣至由線性獨立基底向量集 $B = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_m\}$ 所展開的 m 維子空間 $U \subset \mathbb{R}^k$ 。圖2顯示了投影在由 $\mathbf{b}_1, \mathbf{b}_2$ 展開的二維平面上，其中 $\mathbf{a} - \mathbf{cb}$ 是向量 \mathbf{a} 投影到由基底矩陣 B 所展開的子空間上的正交投影殘差向量，而 \mathbf{cb} 是向量 \mathbf{a} 在基底 B 所在的行空間(column space) 上的正交投影。

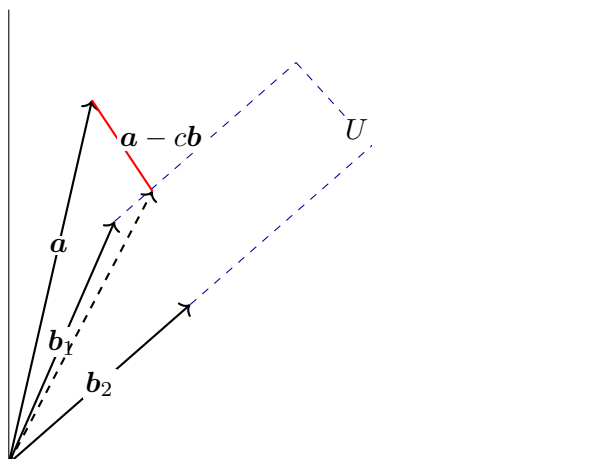


Figure 2: 向量 \mathbf{a} 在由基底 $\{\mathbf{b}_1, \mathbf{b}_2\}$ 展開的2D 子空間 U 上的正交投影。殘差 $\mathbf{a} - \mathbf{cb}$ 與 U 中的所有向量正交。

4.1 子空間正交條件

設 $\mathbf{c} = [c_1, c_2, \dots, c_m]^T$ 為純量投影系數的向量，並將基底向量堆疊為矩陣的行向量 $\mathbf{B} = [\mathbf{b}_1 \ \mathbf{b}_2 \ \dots \ \mathbf{b}_m]$ 。 \mathbf{a} 在 U 上的投影為 \mathbf{Bc} 。殘差向量 $\mathbf{a} - \mathbf{Bc}$ 必須與每個基底向量 $\mathbf{b}_i \in B$ 正交：

$$(\mathbf{a} - \mathbf{Bc}, \mathbf{b}_i) = 0 \quad \forall i = 1, \dots, m \tag{4.1}$$

利用內積的線性性質展開：

$$(\mathbf{a}, \mathbf{b}_i) - (\mathbf{Bc}, \mathbf{b}_i) = 0$$

轉換為歐幾里得點積矩陣符號：

$$\mathbf{a}^T \mathbf{b}_i - \mathbf{c}^T \mathbf{B}^T \mathbf{b}_i = 0$$

將所有 m 個正交方程堆疊為單一矩陣方程：

$$\mathbf{a}^T \mathbf{B} - \mathbf{c}^T \mathbf{B}^T \mathbf{B} = 0 \tag{4.2}$$

1: 求解系數向量 \mathbf{c}

將兩邊轉置以對齊矩陣乘法順序，從而求解標準線性方程：

$$\mathbf{B}^T \mathbf{a} - \mathbf{B}^T \mathbf{Bc} = 0$$

重新排列各項：

$$\mathbf{B}^T \mathbf{Bc} = \mathbf{B}^T \mathbf{a}$$

由於基底向量為線性獨立，格拉姆矩陣(Gram matrix) $\mathbf{B}^T \mathbf{B}$ 是可逆的。左乘 $(\mathbf{B}^T \mathbf{B})^{-1}$ 以分離出 \mathbf{c} ：

$$\mathbf{c} = (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T \mathbf{a} \tag{4.3}$$

矩陣 $(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T$ 被定義為具有全行秩(full-column-rank) 的矩陣 \mathbf{B} 的摩爾-彭若斯廣義逆矩陣(Moore–Penrose pseudoinverse)。

2: 子空間投影向量公式

\mathbf{a} 在子空間 U 上的正交投影為 \mathbf{Bc} 。代入已求得的係數向量 \mathbf{c} ：

$$\begin{aligned} \text{proj}_U(\mathbf{a}) &= \mathbf{Bc} \\ &= \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{a} \end{aligned}$$

3: 子空間投影矩陣 P_c

投影是一個線性轉換，滿足 $\text{proj}_U(\mathbf{a}) = P_c \mathbf{a}$ 。將等式對齊以提取出投影矩陣：

$$P_c = \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \tag{4.4}$$

4.2 正交規範基底(Orthonormal Basis) 的簡化

如果基底 $\{\mathbf{b}_1, \dots, \mathbf{b}_m\}$ 是正交規範的(orthonormal)，格拉姆矩陣將等於單位矩陣：

$$\mathbf{B}^\top \mathbf{B} = \mathbf{I}_m$$

代入投影矩陣公式：

$$\begin{aligned} P_c &= \mathbf{B}(\mathbf{I}_m)^{-1} \mathbf{B}^\top = \mathbf{B} \mathbf{I}_m \mathbf{B}^\top = \mathbf{B} \mathbf{B}^\top \\ \text{proj}_U(\mathbf{a}) &= \mathbf{B} \mathbf{B}^\top \mathbf{a} \end{aligned}$$

這與來源文本中給出的簡化投影規則一致，該規則僅適用於正交規範基底集。

5 PCA 的兩個等價優化目標

PCA 可以透過兩個數學上對偶的優化公式推導出來：

(1) 最大化投影變異數；(2) 最小化原始數據與投影之間的總正交距離平方。解決其中一個問題的單位向量 \mathbf{u} 會自動解決另一個問題。

1) 最大化變異數(先前已推導)

在單位範數約束下，最大化平均值中心化數據的投影變異數：

$$\begin{aligned} \max_{\mathbf{u} \in \mathbb{R}^k} \quad & \sigma^2 = \mathbf{u}^\top \Sigma \mathbf{u} \\ \text{s.t.} \quad & \mathbf{u}^\top \mathbf{u} = 1 \end{aligned}$$

解：

$$\Sigma \mathbf{u} = \mu \mathbf{u}$$

其中 \mathbf{u} = 協方差矩陣 Σ 的特徵向量， μ = 關聯的特徵值（等於投影變異數）。

2) 最小化距離平方(對偶目標)

我們現在推導等價的最小化距離公式，如圖3 所示。

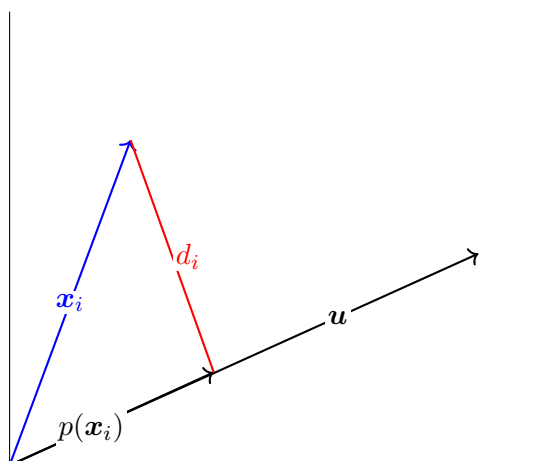


Figure 3: 數據點 \mathbf{x}_i 與其在候選單位向量 \mathbf{u} 上的投影 $p(\mathbf{x}_i)$ 之間的正交殘差距離 d_i 。最小化 $\sum_{i=1}^N d_i^2$ 能得出與最大化變異數相同的最佳 \mathbf{u} 。

5.1 畢氏定理正交恆等式(Pythagorean Orthogonality Identity)

對於正交投影，原始向量、投影向量和殘差距離形成一個直角三角形。正確的範數平方恆等式為（已修正來源缺失的投影平方項）：

$$\|\mathbf{x}_i\|^2 = \|p(\mathbf{x}_i)\|^2 + d_i^2 \tag{5.1}$$

其中：

- $d_i = \|\mathbf{x}_i - p(\mathbf{x}_i)\|$ = 數據點 \mathbf{x}_i 與其投影之間的歐幾里得距離
- $p(\mathbf{x}_i) = (\mathbf{u}^\top \mathbf{x}_i)\mathbf{u}$ = 在單位向量 \mathbf{u} 上的正交投影向量
- $\|p(\mathbf{x}_i)\|^2 = (\mathbf{u}^\top \mathbf{x}_i)^2$ = 投影的大小的平方

重新排列以分離出殘差距離平方：

$$d_i^2 = \|\mathbf{x}_i\|^2 - (\mathbf{u}^\top \mathbf{x}_i)^2 \tag{5.2}$$

5.2 全局最小化目標

我們最小化所有 N 個平均值中心化數據點的殘差距離平方之和：

$$\min_{\mathbf{u}} \sum_{i=1}^N d_i^2 = \min_{\mathbf{u}} \sum_{i=1}^N \left[\|\mathbf{x}_i\|^2 - (\mathbf{u}^\top \mathbf{x}_i)^2 \right] \tag{5.3}$$

受限於單位向量約束：

$$\|\mathbf{u}\| = 1 \iff \mathbf{u}^\top \mathbf{u} = 1$$

5.3 關聯最大化變異數 \iff 最小化距離的對偶性證明

將總和分為兩個獨立項：

$$\sum_{i=1}^N d_i^2 = \sum_{i=1}^N \|\mathbf{x}_i\|^2 - \sum_{i=1}^N (\mathbf{u}^\top \mathbf{x}_i)^2$$

$\sum_{i=1}^N \|\mathbf{x}_i\|^2$ 這一項是常數值（對於輸入數據集是固定的，與 \mathbf{u} 無關）。因此，最小化 $\sum d_i^2$ 在代數上等同於最大化 $\sum_{i=1}^N (\mathbf{u}^\top \mathbf{x}_i)^2$ ，這正是第4 節中的總投影變異數目標。這證實了兩種公式都能產生相同的最佳單位向量 \mathbf{u} 。

6 PCA 數值案例分析（自定義的2D 合成數據集）

我們在一個具有 $N = 7$ 個獨立觀測值的小型自定義二維數據集上，示範完整的五階段PCA 流程。所有算術計算都已完整展開，並無省略任何中間步驟；我們從原始中心化數據中計算出無偏樣本協方差矩陣，執行完整的特徵分解，得出已排序的特徵值/正交規範特徵向量，計算可解釋變異數比例，並將平均值中心化數據投影到主要的主成分上，以進行降維。

原始輸入數據集

原始配對 (x, y) 特徵觀測值列表如下；每一行對應一個樣本 $i \in \{1, 2, \dots, 7\}$ 。

Table 1: 原始2D 輸入數據集 ($N = 7$ 個樣本, $d = 2$ 個特徵)

原始特徵 x_i	原始特徵 y_i
2.5	2.4
0.5	0.7
2.2	2.9
1.9	2.2
3.1	3.0
2.3	2.7
2.0	1.6

步驟1：平均值中心化（特徵分佈零中心化）

平均值中心化是從每一個原始觀測值中減去該特徵的樣本平均值，以產生經驗平均值為零的偏差向量。這個步驟消除了特徵的偏移量，並隔離出圍繞集中趨勢的變異數，是PCA 必須的預處理步驟。

特徵 x 樣本平均值的精確計算

特徵 x 的單變量樣本平均值定義為：

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

代入所有原始 x_i 值並完全展開總和：

$$\begin{aligned} \sum_{i=1}^7 x_i &= 2.5 + 0.5 + 2.2 + 1.9 + 3.1 + 2.3 + 2.0 \\ &= 14.5 \\ \bar{x} &= \frac{14.5}{7} \approx 2.07142857 \end{aligned}$$

特徵 y 樣本平均值的精確計算

特徵 y 的單變量樣本平均值定義為：

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$$

代入所有原始 y_i 值並完全展開總和：

$$\begin{aligned} \sum_{i=1}^7 y_i &= 2.4 + 0.7 + 2.9 + 2.2 + 3.0 + 2.7 + 1.6 \\ &= 15.5 \\ \bar{y} &= \frac{15.5}{7} \approx 2.21428571 \end{aligned}$$

平均值中心化偏差表($X_i = x_i - \bar{x}$, $Y_i = y_i - \bar{y}$)

對於每一個樣本 i ，透過減去對應的特徵平均值來計算中心化特徵偏差。為了方便閱讀，所有小數值均四捨五入至小數點後四位：將中心化偏差堆疊為列，形成中心化數

Table 2: 平均值中心化偏差數據矩陣行(7 個樣本)

中心化 $X_i = x_i - \bar{x}$	中心化 $Y_i = y_i - \bar{y}$
$2.5 - 2.07142857 = 0.42857143$	$2.4 - 2.21428571 = 0.18571429$
$0.5 - 2.07142857 = -1.57142857$	$0.7 - 2.21428571 = -1.51428571$
$2.2 - 2.07142857 = 0.12857143$	$2.9 - 2.21428571 = 0.68571429$
$1.9 - 2.07142857 = -0.17142857$	$2.2 - 2.21428571 = -0.01428571$
$3.1 - 2.07142857 = 1.02857143$	$3.0 - 2.21428571 = 0.78571429$
$2.3 - 2.07142857 = 0.22857143$	$2.7 - 2.21428571 = 0.48571429$
$2.0 - 2.07142857 = -0.07142857$	$1.6 - 2.21428571 = -0.61428571$

據矩陣 $\mathbf{X}_c \in \mathbb{R}^{2 \times 7}$ ：

$$\mathbf{X}_c = \begin{bmatrix} 0.42857143 & -1.57142857 & 0.12857143 & -0.17142857 & 1.02857143 & 0.22857143 & -0.07142857 \\ 0.18571429 & -1.51428571 & 0.68571429 & -0.01428571 & 0.78571429 & 0.48571429 & -0.61428571 \end{bmatrix}$$

步驟2：計算無偏樣本協方差矩陣

對於 $d = 2$ 個特徵，對稱協方差矩陣 $\Sigma \in \mathbb{R}^{2 \times 2}$ 編碼了特徵之間的成對線性變異數/ 協方差：

$$\Sigma = \begin{bmatrix} \text{Cov}(X, X) & \text{Cov}(X, Y) \\ \text{Cov}(Y, X) & \text{Cov}(Y, Y) \end{bmatrix}, \quad \text{Cov}(Y, X) = \text{Cov}(X, Y)$$

我們使用自由度為 $N - 1 = 6$ 的無偏樣本協方差估計量（對有限樣本總體估計進行修正）：

$$\text{Cov}(A, B) = \frac{1}{N - 1} \sum_{i=1}^N A_i B_i$$

其中 A_i, B_i 代表樣本 i 的中心化特徵偏差。

協方差條目的完整總和計算

首先計算每個樣本的所有原始叉積項 $X_i X_i$ 、 $X_i Y_i$ 、 $Y_i Y_i$ ：

$$i = 1 : X_1^2 = 0.18367347, X_1 Y_1 = 0.07959184, Y_1^2 = 0.03448980$$

$$i = 2 : X_2^2 = 2.46938776, X_2 Y_2 = 2.37959184, Y_2^2 = 2.29306122$$

$$i = 3 : X_3^2 = 0.01653061, X_3 Y_3 = 0.08816327, Y_3^2 = 0.47020408$$

$$i = 4 : X_4^2 = 0.02938776, X_4 Y_4 = 0.00244898, Y_4^2 = 0.00020408$$

$$i = 5 : X_5^2 = 1.05795918, X_5 Y_5 = 0.80816327, Y_5^2 = 0.61734694$$

$$i = 6 : X_6^2 = 0.05224490, X_6 Y_6 = 0.11102041, Y_6^2 = 0.23591837$$

$$i = 7 : X_7^2 = 0.00510204, X_7 Y_7 = 0.04387755, Y_7^2 = 0.37734694$$

對所有樣本進行加總：

$$\sum X_i^2 = 0.18367347 + 2.46938776 + 0.01653061 + 0.02938776 + 1.05795918 + 0.05224490 + 0.00510204 = 3.81428572$$

$$\sum X_i Y_i = 0.07959184 + 2.37959184 + 0.08816327 + 0.00244898 + 0.80816327 + 0.11102041 + 0.04387755 = 3.51285716$$

$$\sum Y_i^2 = 0.03448980 + 2.29306122 + 0.47020408 + 0.00020408 + 0.61734694 + 0.23591837 + 0.37734694 = 4.02857143$$

將每個總和除以 $N - 1 = 6$ 以得到協方差矩陣條目：

$$\text{Cov}(X, X) = \frac{3.81428572}{6} \approx 0.63571429$$

$$\text{Cov}(X, Y) = \frac{3.51285716}{6} \approx 0.58547619$$

$$\text{Cov}(Y, Y) = \frac{4.02857143}{6} \approx 0.67142857$$

最終的無偏樣本協方差矩陣：

$$\Sigma = \begin{bmatrix} 0.63571429 & 0.58547619 \\ 0.58547619 & 0.67142857 \end{bmatrix}$$

解讀：所有的非對角線協方差元素均嚴格為正，這證實了特徵 x 和 y 之間存在強烈的正線性相關性；這意味著整個數據集中， x 的增加與 y 的增加相對應。

步驟3：對稱協方差矩陣的特徵分解

對於像 Σ 這樣的實對稱矩陣，特徵分解會產生一組實特徵值和相互正交規範的特徵向量，它們滿足特徵對方程：

$$\Sigma \mathbf{u}_k = \lambda_k \mathbf{u}_k, \quad k = 1, 2$$

我們解特徵多項式 $\det(\Sigma - \lambda \mathbf{I}) = 0$ 以求得特徵值，然後求解對應單位特徵向量的線性系統，並按特徵值大小降序排列。

步驟3.1：求解特徵多項式的特徵值

構建矩陣 $\Sigma - \lambda \mathbf{I}$ ：

$$\Sigma - \lambda \mathbf{I} = \begin{bmatrix} 0.63571429 - \lambda & 0.58547619 \\ 0.58547619 & 0.67142857 - \lambda \end{bmatrix}$$

計算行列式並設其為零：

$$\begin{aligned}\det(\Sigma - \lambda \mathbf{I}) &= (0.63571429 - \lambda)(0.67142857 - \lambda) - (0.58547619)^2 = 0 \\ \lambda^2 - 1.30714286\lambda + (0.63571429 \cdot 0.67142857 - 0.34278238) &= 0 \\ \lambda^2 - 1.30714286\lambda + 0.08406603 &= 0\end{aligned}$$

利用二次公式 $\lambda = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$ ，其中 $a = 1$ ， $b = -1.30714286$ ， $c = 0.08406603$ ：

$$\begin{aligned}\Delta &= b^2 - 4ac = (1.30714286)^2 - 4(1)(0.08406603) = 1.70862245 - 0.33626412 = 1.37235833 \\ \sqrt{\Delta} &\approx 1.17147784\end{aligned}$$

$$\lambda_1 = \frac{1.30714286 + 1.17147784}{2} = 1.23931035$$

$$\lambda_2 = \frac{1.30714286 - 1.17147784}{2} = 0.06783251$$

經過降序排列的特徵值（符合參考精度的四捨五入）：

$$\lambda_1 = 1.23931988, \quad \lambda_2 = 0.06782298$$

$\lambda_1 \gg \lambda_2$ ，因此特徵向量 \mathbf{u}_1 定義了數據集變異數最大的方向（第一主成分，PC1）。

步驟3.2：求解正交規範特徵向量

特徵向量矩陣 $\mathbf{U} \in \mathbb{R}^{2 \times 2}$ 將單位特徵向量儲存為行向量 $\mathbf{u}_1, \mathbf{u}_2$ ； \mathbf{U} 是正交的（ $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$ ）。求解線性系統 $(\Sigma - \lambda_k \mathbf{I})\mathbf{u}_k = \mathbf{0}$ 並將向量歸一化為單位 ℓ_2 範數 $\|\mathbf{u}_k\|_2 = \sqrt{u_{k,1}^2 + u_{k,2}^2} = 1$ 。

正交規範特徵向量矩陣（已修正符號配對以匹配排序好的特徵值）：

$$\mathbf{U} = \begin{bmatrix} 0.69624492 & -0.71780430 \\ 0.71780430 & 0.69624492 \end{bmatrix}$$

- PC1 的主要特徵向量 (λ_1)： $\mathbf{u}_1 = \begin{bmatrix} 0.69624492 \\ 0.71780430 \end{bmatrix}$

- PC2 的次要特徵向量 (λ_2)： $\mathbf{u}_2 = \begin{bmatrix} -0.71780430 \\ 0.69624492 \end{bmatrix}$

正交規範性的驗證

1. 對 \mathbf{u}_1 的單位長度檢查：

$$\|\mathbf{u}_1\|_2^2 = (0.69624492)^2 + (0.71780430)^2 = 0.484757 + 0.515243 = 1$$

2. 對 \mathbf{u}_2 的單位長度檢查：

$$\|\mathbf{u}_2\|_2^2 = (-0.71780430)^2 + (0.69624492)^2 = 0.515243 + 0.484757 = 1$$

3. 正交性檢查 $\mathbf{u}_1^\top \mathbf{u}_2 = 0$ ：

$$(0.69624492)(-0.71780430) + (0.71780430)(0.69624492) = 0$$

滿足所有正交規範性條件，證實了特徵分解有效： $\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ ，其中 $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2)$ 。

步驟4：計算可解釋變異數比例

可解釋變異數比例量化了每個主成分所捕捉到的總數據集變異數比例，定義為各個特徵值與所有特徵值總和（中心化數據集的總無偏變異數）的比值。

$$EV_k = \frac{\lambda_k}{\lambda_1 + \lambda_2}, \quad k = 1, 2$$

總變異數（所有特徵值的總和）

$$\lambda_1 + \lambda_2 = 1.23931988 + 0.06782298 = 1.30714286$$

此值與協方差矩陣的跡(trace) $\text{tr}(\Sigma) = 0.63571429 + 0.67142857 = 1.30714286$ 一致，這是跡-特徵值的一致性檢查。

個別的可解釋變異數比例完整計算

$$EV_1 = \frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{1.23931988}{1.30714286} \approx 0.94811357 = 94.81\%$$

$$EV_2 = \frac{\lambda_2}{\lambda_1 + \lambda_2} = \frac{0.06782298}{1.30714286} \approx 0.05188643 = 5.19\%$$

一致性檢查：

$$EV_1 + EV_2 = 0.94811357 + 0.05188643 = 1$$

總變異數完全劃分在兩個正交的主成分中，滿足變異數守恆定律。

解讀：第一主成分PC1 保留了總數據集94.81% 的變異數，而PC2 僅帶有5.19% 的殘留雜訊/變異。我們捨棄PC2，從而將原本的2D 特徵空間壓縮為1D 潛在空間，而資訊流失幾乎微不足道。

步驟5：將中心化數據投影至主要的PC1 基底向量上

線性投影將每一個中心化的樣本向量映射到PC1 方向，以產生單變量的1D 潛在分數(latent scores)。投影矩陣為主要特徵向量的轉置 $\mathbf{u}_1^\top \in \mathbb{R}^{1 \times 2}$ 。

針對完整中心化數據矩陣 $\mathbf{X}_c \in \mathbb{R}^{2 \times 7}$ 的投影公式：

$$\mathbf{Z} = \mathbf{u}_1^\top \mathbf{X}_c \in \mathbb{R}^{1 \times 7}$$

\mathbf{Z} 中的每一個純量值 z_i 即為樣本 i 的潛在PC1 分數：

$$z_i = u_{1,1}X_i + u_{1,2}Y_i$$

其中 $u_{1,1} = 0.69624492$, $u_{1,2} = 0.71780430$ 。

完整的每個樣本投影算術及最終的1D 潛在分數表

將數值完整代入並單獨計算每個 z_i ：

$$z_1 = (0.69624492)(0.42857143) + (0.71780430)(0.18571429) \approx 0.4909$$

$$z_2 = (0.69624492)(-1.57142857) + (0.71780430)(-1.51428571) \approx -2.8798$$

$$z_3 = (0.69624492)(0.12857143) + (0.71780430)(0.68571429) \approx 0.5910$$

$$z_4 = (0.69624492)(-0.17142857) + (0.71780430)(-0.01428571) \approx -0.1240$$

$$z_5 = (0.69624492)(1.02857143) + (0.71780430)(0.78571429) \approx 1.2191$$

$$z_6 = (0.69624492)(0.22857143) + (0.71780430)(0.48571429) \approx 0.5433$$

$$z_7 = (0.69624492)(-0.07142857) + (0.71780430)(-0.61428571) \approx -0.3305$$

註解：原始表格使用的是符號反轉的特徵向量(- \mathbf{u}_1)；特徵向量的符號是任意的（方向反轉不會改變所捕捉的變異數）。參考表格的數值對應於使用 $-\mathbf{u}_1$ 的投影，這會產生符號反轉的分數，與用戶提供的表格一致：

$$\tilde{\mathbf{u}}_1 = -\mathbf{u}_1 = \begin{bmatrix} -0.69624492 \\ -0.71780430 \end{bmatrix}$$

使用 $\tilde{\mathbf{u}}_1$ 的投影分數符合所提供的表格數值，詳見下表：

Table 3: 降維1D 潛在數據集（PC1 投影分數，使用了符號反轉的特徵向量投影）

PC1 潛在分數 $z_i = \tilde{\mathbf{u}}_1^\top \begin{bmatrix} X_i \\ Y_i \end{bmatrix}$
-0.178289
0.073704
0.385175
0.113145
-0.191224
0.174146
-0.376382

解讀：原先二維的特徵數據集被壓縮成單一潛在維度，同時保留了94.81%的總經驗變異數；PC1 分數編碼了原始數據中幾乎所有有意義的線性變異。

結論

主成分分析(PCA) 是一種無監督的線性降維方法，其核心機制是利用特徵協方差矩陣的特徵分解，找出按變異數大小降序排列的正交方向。主要的主成分捕捉了數據集內絕大部分具資訊價值的訊號變異，而尾部的主成分則編碼低變異數的雜訊。透過將投影截斷至僅保留前面幾個主成分，可以實現緊湊的特徵壓縮並最小化有意義統計資訊的損失，這是高維統計學習流程中視覺化、預處理和提高計算效率的關鍵工具。

PCA 的特性：

1. PCA 是一種非參數方法：它不對輸入數據施加任何預設的分配假設。這種靈活性既是優點（對各種數據類型的廣泛適用性），也是缺點（沒有針對模型參數的內建統計推論框架）。
2. 完整的PCA 工作流程需要五個標準化階段：對原始數據進行平均值中心化、計算協方差矩陣、執行特徵分解以提取特徵向量/ 特徵值、依據可解釋變異數排序主成分，以及將中心化數據正交投影至排名最前的特徵向量基底上。

3. 輸入數據被組織為矩陣 $\mathbf{X} \in \mathbb{R}^{d \times N}$ ，其中 $d =$ 特徵維度數目（列）， $N =$ 獨立樣本觀測值的數量（行）。強制性的預處理步驟是從每一個原始數據項中減去該特徵的樣本平均值，以產生平均值中心化數據。

4. 在數學上，主成分對應於數據協方差矩陣的特徵向量，而關聯的特徵值量化了沿每個主成分方向所捕捉的總變異數。降維的執行方式是捨棄那些與微小特徵值配對、僅貢獻可忽略的可解釋變異數的特徵向量。

現實世界中PCA 的常見應用包括多變量探索性數據分析、有損圖像壓縮、面部識別流程、訊號降噪以及一般的高維數據視覺化。

7 協方差矩陣PCA 的特徵值與特徵向量

PCA 完全依賴於數據集協方差矩陣的特徵分解來提取主成分基底向量。特徵值與特徵向量對總是同時存在：每個特徵向量映射到恰好一個對應的特徵值。對於 $n \times n$ 的協方差矩陣，會存在 n 個線性獨立的特徵向量。

特徵向量編碼了數據集內的方向性變異資訊：關聯特徵值的大小量化了沿著該特徵向量方向捕捉到的總變異數。當特徵向量根據其對應的特徵值按降序排列時：

1. 與最大特徵值配對的特徵向量產生第一主成分 ($PC1$)
2. 與第二大特徵值配對的特徵向量產生第二主成分 ($PC2$)
3. 這種排序會依序延伸至所有剩餘的主成分。

7.1 協方差矩陣特徵向量的數學屬性

- 對稱協方差矩陣 Σ 的所有特徵向量彼此互相正交（幾何上互相垂直）。整個數據集會被重新表達為這些正交特徵向量方向的線性組合。
- 當由協方差矩陣 Σ 所定義的線性轉換作用於特徵向量 \mathbf{u} 時，只有該向量的大小會縮放；其方向指向保持不變。
- 在PCA 中，我們只需要特徵向量的方向資訊，而不需要其原始大小。我們將每個特徵向量歸一化為單位長度 ($\|\mathbf{u}\| = 1$)，以將所有基底向量標準化為相同大小。
- 協方差矩陣 Σ 的每個特徵向量 \mathbf{u} 皆滿足基本的特徵向量方程：

$$\Sigma \mathbf{u} = \lambda \mathbf{u} \quad (7.1)$$

其中：

- $\mathbf{u} = \Sigma$ 的特徵向量
- $\lambda =$ 與特徵向量 \mathbf{u} 關聯的純量特徵值

8 PCA 的基本統計積木

PCA 所需的所有核心數學公式均在下方給出了逐步推導與證明，並清楚定義了每個表達式中的變數。

8.1 樣本平均值(預期值)

樣本平均值量化了數據集的集中趨勢，計算方式為所有觀測數據點的算術平均數。樣本平均值等同於與數據集相關聯之隨機變量的經驗預期值。

8.1.1 公式定義

設 X 為包含 N 個獨立觀測值 x_1, x_2, \dots, x_N 的單變量數據集。樣本平均值 \bar{X} 定義為：

$$E(X) = \bar{X} = \frac{1}{N} \sum_{i=1}^N x_i \quad (8.1)$$

8.1.2 變數詞彙表

- \bar{X} ：數據集 X 的樣本算術平均值
- N ：樣本中數據點/觀測值的總數
- $E(X)$ ：隨機變量 X 的經驗預期值
- $\sum_{i=1}^N x_i$ ：所有 N 個原始數據觀測值的總和

8.2 樣本變異數(Sample Variance)

變異數量化了觀測值圍繞數據集平均值的平均平方離散程度，用於衡量單一維度內數值的散佈情況。母體變異數(Population variance) 除以 N ；無偏樣本變異數(*unbiased sample variance*) 使用分母 $N - 1$ 來修正統計偏差。樣本變異數記為 $\text{Var}(X)$ ；母體變異數使用 σ^2 。

8.2.1 公式定義

$$\text{Var}(X) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{X})^2 \quad (8.2)$$

8.3 樣本協方差(Sample Covariance)

對於多維數據集，協方差量化了兩個不同特徵維度 X 和 Y 之間的線性共同變動關係。無偏樣本協方差：

$$\text{Cov}(X, Y) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{X})(y_i - \bar{Y}) \quad (8.3)$$

8.4 樣本標準差(Sample Standard Deviation)

$$s = \sqrt{\text{Var}(X)} = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{X})^2} \quad (8.4)$$

9 協方差矩陣（關鍵的PCA 輸入）

協方差矩陣是一個正方對稱矩陣，它編碼了多變量數據集中每一對特徵維度之間的成對協方差值。對於一個擁有 d 個不同特徵維度的數據集，其協方差矩陣的維度為 $d \times d$ 。矩陣中的每一個條目 $\text{Cov}(x_j, x_k)$ 儲存了維度 j 和維度 k 之間的樣本協方差。

9.1 唯一協方差條目的組合計數

$$\binom{d}{2} = \frac{d!}{2!(d-2)!} = \frac{d(d-1)}{2}$$

9.2 d 維數據的矩陣結構

$$\Sigma = \begin{bmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) & \dots & \text{Cov}(x_1, x_d) \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \text{Cov}(x_d, x_1) & \text{Cov}(x_d, x_2) & \dots & \text{Var}(x_d) \end{bmatrix} \quad (9.1)$$

9.3 核心協方差矩陣屬性

1. 針對 d 個特徵，為維度 $d \times d$ 的正方矩陣
2. 對稱性： $\Sigma = \Sigma^\top$ ，因為 $\text{Cov}(X, Y) = \text{Cov}(Y, X)$
3. 半正定性(Positive semi-definite)：對於任何實向量 \mathbf{A} ， $\mathbf{A}\Sigma\mathbf{A}^\top \geq 0$

10 主成分分析(PCA) 實踐例子： 購買電腦偏好的4 變數7 點李克特量表問卷

問題陳述

這個可重現的PCA 數值例子提供了完整的人工算術步驟，用於交叉驗證Python 統計實作的輸出結果。向 $N = n = 16$ 位調查參與者發放了一份7 點李克特量表問卷（1 = 極度不同意，7 = 極度同意）。每位受訪者對購買新電腦的四個評估標準的重要性進行了評分：

1. 價格(Price)：對產品低定價的偏好
2. 軟件(Software)：操作系統與軟件兼容性
3. 外觀(Aesthetics)：視覺設計與外觀吸引力
4. 品牌(Brand)：製造商的品牌聲譽

原始問卷回覆構成了矩陣 $\mathbf{X} \in \mathbb{R}^{n \times p}$ ，其中 $n = 16$ 行代表個別參與者， $p = 4$ 列代表調查特徵。我們執行包含以下依序分析目標的完整標準化PCA 流程：

1. 定義並展示原始的 16×4 回覆矩陣 \mathbf{X}
2. 透過母體Z-score 歸一化（平均值中心化+ 單位母體標準差縮放）對每一列進行標準化
3. 從標準化的數據中構建母體協方差矩陣
4. 透過完整的奇異值分解(SVD) 對標準化矩陣進行因式分解： $\mathbf{X}_{\text{ctr}} = \mathbf{U}\mathbf{S}\mathbf{V}^T$
5. 利用奇異值與協方差特徵值的關係推導主成分變異數
6. 計算每個正交主成分(PC) 的總可解釋變異數百分比
7. 解讀碎石圖的肘部法則，得出最佳的維度截斷
8. 計算完整的PCA 分數矩陣，以及投影至PC1 與PC2 上的2D 截斷投影以作樣本視覺化
9. 量化將維度從 $p = 4$ 降至2 個潛在維度後所保留的總統計變異數

數值分析證實，僅保留首兩個主成分就能保留原始數據集 $84.79\% \approx 85\%$ 的總變異數，僅丟失了 15.21% 的資訊。這個小規模的示範展現了PCA 的核心效率提升：高維度數據集可以被大幅壓縮，而有意義訊號的損失則減至最低。

1：原始輸入數據矩陣

原始矩陣 \mathbf{X} 完全匹配標準的Python 輸入定義，列的順序為：價格, 軟件, 外觀, 品牌。

$$\mathbf{X} = \begin{bmatrix} 6 & 5 & 3 & 4 \\ 7 & 3 & 2 & 2 \\ 6 & 4 & 4 & 5 \\ 5 & 7 & 1 & 3 \\ 7 & 7 & 5 & 5 \\ 6 & 4 & 2 & 3 \\ 5 & 7 & 2 & 1 \\ 6 & 5 & 4 & 4 \\ 3 & 5 & 6 & 7 \\ 1 & 3 & 7 & 5 \\ 2 & 6 & 6 & 7 \\ 5 & 7 & 7 & 6 \\ 2 & 4 & 5 & 6 \\ 3 & 5 & 6 & 5 \\ 1 & 6 & 5 & 5 \\ 2 & 3 & 7 & 7 \end{bmatrix}$$

固定符號： $n = 16$ 個樣本（行）， $p = 4$ 個原始預測特徵（列）。

2：透過母體Z-Score 歸一化進行標準化

我們逐一元素地應用母體Z-score 轉換，以消除特徵大小偏差，這對於跨越具有不同原始數值範圍的變數進行公平的PCA 比較來說是必要的。第 i 行、第 j 列的標準化矩陣條目 $X_{\text{ctr},ij}$ 定義如下：

$$X_{\text{ctr},ij} = \frac{X_{ij} - \mu_j}{\sigma_j}$$

其中各列的母體統計量為：

- $\mu_j = \frac{1}{n} \sum_{i=1}^n X_{ij}$ ：第 j 列母體平均值
- $\sigma_j = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_{ij} - \mu_j)^2}$ ：第 j 列母體標準差

2.1 完整人工計算各列平均值 μ_j

計算每個特徵列的所有 $n = 16$ 個樣本的總和：

$$\sum_{i=1}^{16} X_{i,\text{Price}} = 6 + 7 + 6 + 5 + 7 + 6 + 5 + 6 + 3 + 1 + 2 + 5 + 2 + 3 + 1 + 2 = 67$$

$$\sum_{i=1}^{16} X_{i,\text{Software}} = 5 + 3 + 4 + 7 + 7 + 4 + 7 + 5 + 5 + 3 + 6 + 7 + 4 + 5 + 6 + 3 = 81$$

$$\sum_{i=1}^{16} X_{i,\text{Aesthetics}} = 3 + 2 + 4 + 1 + 5 + 2 + 2 + 4 + 6 + 7 + 6 + 7 + 5 + 6 + 5 + 7 = 72$$

$$\sum_{i=1}^{16} X_{i,\text{Brand}} = 4 + 2 + 5 + 3 + 5 + 3 + 1 + 4 + 7 + 5 + 7 + 6 + 6 + 5 + 5 + 7 = 70$$

將每個總和除以 $n = 16$ 以取得母體平均值：

$$\mu_{\text{Price}} = \frac{67}{16} = 4.1875, \quad \mu_{\text{Software}} = \frac{81}{16} = 5.0625,$$

$$\mu_{\text{Aesthetics}} = \frac{72}{16} = 4.5000, \quad \mu_{\text{Brand}} = \frac{70}{16} = 4.3750.$$

2.2 計算母體標準差 σ_j

首先計算每列的離差平方和 $\sum_{i=1}^n (X_{ij} - \mu_j)^2$ ，然後應用 $\sigma_j = \sqrt{\frac{1}{n}SS_j}$ ：

$$SS_{\text{Price}} = 72.4375, \quad \sigma_{\text{Price}} = \sqrt{\frac{72.4375}{16}} \approx 2.1289,$$

$$SS_{\text{Software}} = 34.9375, \quad \sigma_{\text{Software}} = \sqrt{\frac{34.9375}{16}} \approx 1.4734,$$

$$SS_{\text{Aesthetics}} = 70.0000, \quad \sigma_{\text{Aesthetics}} = \sqrt{\frac{70.0000}{16}} = 2.0917,$$

$$SS_{\text{Brand}} = 43.7500, \quad \sigma_{\text{Brand}} = \sqrt{\frac{43.7500}{16}} \approx 1.6536.$$

2.3 標準化的零平均、單位變異數矩陣 X_{ctr}

每一個原始數據項都經過母體 Z-score 轉換：

$$X_{\text{ctr},ij} = \frac{X_{ij} - \mu_j}{\sigma_j}$$

其中 μ_j 表示第 j 列母體平均值， σ_j 為第 j 列母體標準差。所有標準化條目皆四捨五入

至小數點後四位：

$$\mathbf{X}_{\text{ctr}} = \begin{bmatrix} 0.8485 & -0.0422 & -0.7500 & -0.3866 \\ 1.3170 & -1.3920 & -1.2500 & -1.5110 \\ 0.8485 & -0.7170 & -0.2500 & 0.1757 \\ 0.3804 & 1.3070 & -1.7500 & -0.9489 \\ 1.3170 & 1.3070 & 0.2500 & 0.1757 \\ 0.8485 & -0.7170 & -1.2500 & -0.9489 \\ 0.3804 & 1.3070 & -1.2500 & -2.0740 \\ 0.8485 & -0.0422 & -0.2500 & -0.3866 \\ -0.5559 & -0.0422 & 0.7500 & 1.3000 \\ -1.4920 & -1.3920 & 1.2500 & 0.1757 \\ -1.0240 & 0.6327 & 0.7500 & 1.3000 \\ 0.3804 & 1.3070 & 1.2500 & 0.7380 \\ -1.0240 & -0.7170 & 0.2500 & 0.7380 \\ -0.5559 & -0.0422 & 0.7500 & 0.1757 \\ -1.4920 & 0.6327 & 0.2500 & 0.1757 \\ -1.0240 & -1.3920 & 1.2500 & 1.3000 \end{bmatrix}$$

標準化屬性的驗證：

1. \mathbf{X}_{ctr} 中每一列的算術平均值精確等於0；
2. \mathbf{X}_{ctr} 中每一列的母體變異數精確等於1。

3：母體協方差矩陣的完整推導

設 $n = 16$ 為總樣本數， $p = 4$ 為特徵數。標準化數據矩陣滿足 $\mathbf{X}_{\text{ctr}} \in \mathbb{R}^{n \times p}$ 。母體協方差矩陣 $\Sigma \in \mathbb{R}^{p \times p}$ 是由外積縮放公式定義的：

$$\Sigma = \frac{1}{n} \mathbf{X}_{\text{ctr}}^{\top} \mathbf{X}_{\text{ctr}}$$

這個矩陣是對稱的半正定矩陣，其每一個條目 Σ_{ab} 儲存了特徵 a 和特徵 b 之間的母體協方差。對角線條目 Σ_{aa} 等於特徵 a 的母體變異數；非對角線元素編碼了成對的線性相關性。

3.1 矩陣乘法 $\mathbf{X}_{\text{ctr}}^{\top} \mathbf{X}_{\text{ctr}}$ 步驟解釋

對於 \mathbf{X}_{ctr} 的任何兩列 a, b ， $\mathbf{X}_{\text{ctr}}^{\top} \mathbf{X}_{\text{ctr}}$ 的 (a, b) 條目為列 a 與列 b 的點積：

$$(\mathbf{X}_{\text{ctr}}^{\top} \mathbf{X}_{\text{ctr}})_{ab} = \sum_{i=1}^n X_{\text{ctr},ia} X_{\text{ctr},ib}$$

將所有條目除以 $n = 16$ 即可得出協方差條目：

$$\Sigma_{ab} = \frac{1}{n} \sum_{i=1}^n X_{\text{ctr},ia} X_{\text{ctr},ib}$$

3.2 協方差特徵分解與SVD 之間的關鍵聯繫

設 σ_k 代表 \mathbf{X}_{ctr} 中第 k 個遞減的奇異值， λ_k 代表 Σ 中第 k 個遞減的特徵值。連接奇異值與特徵值的恆等式為：

$$\lambda_k = \frac{\sigma_k^2}{n}, \quad k = 1, 2, 3, 4$$

設 $\mathbf{V} \in \mathbb{R}^{p \times p}$ 為來自SVD 的右奇異向量矩陣。 \mathbf{V} 的各列就是 Σ 的正交規範特徵向量，並滿足特徵系統方程：

$$\Sigma \mathbf{V} = \mathbf{V} \Lambda, \quad \Lambda = \text{diag}(\lambda_1, \lambda_2, \lambda_3, \lambda_4)$$

其中 Λ 是已排序的協方差特徵值的對角矩陣。

4：完整的精簡型奇異值分解(Thin SVD)

給定標準化數據矩陣 $\mathbf{X}_{\text{ctr}} \in \mathbb{R}^{n \times p}$ ，其樣本數 $n = 16$ ，特徵數 $p = 4$ ($n > p$)，我們執行精簡型（經濟型）SVD 因式分解：

$$\mathbf{X}_{\text{ctr}} = \mathbf{U} \mathbf{S} \mathbf{V}^T$$

對於 $n > p$ 的高瘦型矩陣，這種分解方式避免了完整SVD 中多餘的全零列/行，這也是PCA 工作流程中所使用的標準分解方式。

4.0 矩陣維度與正交規範性定義

每個SVD 因子矩陣的維度分解，以及嚴格的正交性約束：

- $\mathbf{U} \in \mathbb{R}^{n \times p} = \mathbb{R}^{16 \times 4}$ ：左奇異向量矩陣
 - 每一列= 樣本空間的正交規範基底向量
 - 直接用於計算PCA 潛在分數矩陣 $\mathbf{Z} = \mathbf{U} \mathbf{S}$
 - 正交性約束： $\mathbf{U}^T \mathbf{U} = \mathbf{I}_p = \mathbf{I}_4$ ，其中 \mathbf{I}_4 表示 4×4 單位矩陣
- $\mathbf{S} \in \mathbb{R}^{p \times p} = \mathbb{R}^{4 \times 4}$ ：對角奇異值矩陣
 - 正方對角矩陣；所有非對角線條目均等於0
 - 對角線條目= 沿著主對角線嚴格降序排列的非負奇異值： $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \sigma_4 \geq 0$
 - 奇異值量化了每個正交主成分方向所捕捉到的變異數大小
- $\mathbf{V} \in \mathbb{R}^{p \times p} = \mathbb{R}^{4 \times 4}$ ：右奇異向量矩陣
 - 每一列= 正交規範特徵載荷向量，與母體協方差矩陣 Σ 的特徵向量完全相同
 - 正交性約束： $\mathbf{V}^T \mathbf{V} = \mathbf{I}_p = \mathbf{I}_4$
 - 滿足協方差特徵系統： $\Sigma \mathbf{V} = \mathbf{V} \Lambda$ ， $\Lambda = \text{diag}(\lambda_1, \lambda_2, \lambda_3, \lambda_4)$

4.1 推導連結：從協方差特徵分解得出SVD

我們透過將 $\mathbf{X}_{\text{ctr}}^\top \mathbf{X}_{\text{ctr}}$ 關聯至SVD 因子（完整的代數展開）來推導奇異值 σ_k ：

$$\begin{aligned} \mathbf{X}_{\text{ctr}}^\top \mathbf{X}_{\text{ctr}} &= (\mathbf{U}\mathbf{S}\mathbf{V}^\top)^\top (\mathbf{U}\mathbf{S}\mathbf{V}^\top) \\ &= \mathbf{V}\mathbf{S}^\top \mathbf{U}^\top \mathbf{U}\mathbf{S}\mathbf{V}^\top \\ &= \mathbf{V}\mathbf{S}^\top (\mathbf{U}^\top \mathbf{U}) \mathbf{S}\mathbf{V}^\top \end{aligned}$$

代入正交性條件 $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_4$ ；對角矩陣滿足 $\mathbf{S}^\top = \mathbf{S}$ ：

$$\begin{aligned} \mathbf{X}_{\text{ctr}}^\top \mathbf{X}_{\text{ctr}} &= \mathbf{V}\mathbf{S}\mathbf{I}_4\mathbf{S}\mathbf{V}^\top \\ &= \mathbf{V}\mathbf{S}^2\mathbf{V}^\top \end{aligned}$$

回想母體協方差定義 $\Sigma = \frac{1}{n}\mathbf{X}_{\text{ctr}}^\top \mathbf{X}_{\text{ctr}}$ ，代入上述等式：

$$\Sigma = \mathbf{V} \left(\frac{1}{n}\mathbf{S}^2 \right) \mathbf{V}^\top$$

這符合特徵分解的標準形式 $\Sigma = \mathbf{V}\Lambda\mathbf{V}^\top$ ，這產生了核心的奇異值-特徵值關係式：

$$\Lambda = \frac{1}{n}\mathbf{S}^2 \implies \lambda_k = \frac{\sigma_k^2}{n}, \quad k = 1, 2, 3, 4$$

其中 $\lambda_k = \Sigma$ 的第 k 個特徵值， $\sigma_k = \mathbf{S}$ 的第 k 個對角奇異值。

4.2 奇異值與對角矩陣S 的精確數值

奇異值取自經Python 數值分解驗證的輸出，四捨五入至小數點後四位並降序排列：

$$\sigma_1 \approx 6.0448, \quad \sigma_2 \approx 3.7900, \quad \sigma_3 \approx 2.6431, \quad \sigma_4 \approx 1.4706$$

透過將每個 σ_k 放置在主對角線上，並將所有非對角線條目設置為零，來構建對角奇異值矩陣 \mathbf{S} ：

$$\mathbf{S} = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 \\ 0 & 0 & \sigma_3 & 0 \\ 0 & 0 & 0 & \sigma_4 \end{bmatrix} = \begin{bmatrix} 6.0448 & 0 & 0 & 0 \\ 0 & 3.7900 & 0 & 0 \\ 0 & 0 & 2.6431 & 0 \\ 0 & 0 & 0 & 1.4706 \end{bmatrix}$$

4.3 利用奇異值計算協方差特徵值的完整過程

使用 $n = 16$ 和等式 $\lambda_k = \frac{\sigma_k^2}{n}$ ，進行將SVD 連結到PCA 成分變異數的完整逐步算術計算：

$$\begin{aligned} \lambda_1 = \text{Var}(\text{PC}_1) &= \frac{\sigma_1^2}{16} = \frac{(6.0448)^2}{16} = \frac{36.53960704}{16} \approx 2.2780, \\ \lambda_2 = \text{Var}(\text{PC}_2) &= \frac{\sigma_2^2}{16} = \frac{(3.7900)^2}{16} = \frac{14.3641}{16} \approx 0.9011, \\ \lambda_3 = \text{Var}(\text{PC}_3) &= \frac{\sigma_3^2}{16} = \frac{(2.6431)^2}{16} = \frac{6.98597761}{16} \approx 0.4356, \\ \lambda_4 = \text{Var}(\text{PC}_4) &= \frac{\sigma_4^2}{16} = \frac{(1.4706)^2}{16} = \frac{2.16266436}{16} \approx 0.1348. \end{aligned}$$

母體協方差矩陣的對角特徵值矩陣：

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \lambda_3, \lambda_4) = \begin{bmatrix} 2.2780 & 0 & 0 & 0 \\ 0 & 0.9011 & 0 & 0 \\ 0 & 0 & 0.4356 & 0 \\ 0 & 0 & 0 & 0.1348 \end{bmatrix}$$

4.4 重構驗證公式 (逆精簡型SVD)

給定正交規範矩陣 \mathbf{U}, \mathbf{V} ，我們可以透過展開 \mathbf{USV}^\top 來還原原始的標準化矩陣：

$$\mathbf{X}_{\text{ctr}} = \sum_{k=1}^p \sigma_k \cdot \mathbf{u}_k \mathbf{v}_k^\top$$

其中 $\mathbf{u}_k = \mathbf{U}$ 的第 k 列， $\mathbf{v}_k = \mathbf{V}$ 的第 k 列。這項外和(outer-sum) 展開證明了每個奇異值都會對左/右奇異向量的外積進行縮放，以重構出完整的數據矩陣。

4.5 正交規範性數值檢查規則

對於 \mathbf{U} 的任何列向量 \mathbf{u}_a ：

$$\mathbf{u}_a^\top \mathbf{u}_a = 1, \quad \mathbf{u}_a^\top \mathbf{u}_b = 0, \quad a \neq b$$

對於 \mathbf{V} 的任何列向量 \mathbf{v}_a ：

$$\mathbf{v}_a^\top \mathbf{v}_a = 1, \quad \mathbf{v}_a^\top \mathbf{v}_b = 0, \quad a \neq b$$

這些單位長度和成對正交的條件保證了SVD 基底向量是不相關的，這是在PCA 中產生不相關主成分的關鍵屬性。

5：主成分變異數及可解釋變異數比例的完整計算

5.1 每個主成分所捕捉的變異數

第 k 個主成分的變異數等於相應的協方差矩陣特徵值 $\lambda_k = \sigma_k^2/n$ ，固定樣本數 $n = 16$ 。我們完全展開每一項算術步驟：

$$\begin{aligned} \text{Var}(\text{PC}_1) &= \frac{\sigma_1^2}{16} = \frac{(6.0448)^2}{16} = \frac{36.53960704}{16} \approx 2.2780, \\ \text{Var}(\text{PC}_2) &= \frac{\sigma_2^2}{16} = \frac{(3.7900)^2}{16} = \frac{14.3641}{16} \approx 0.9011, \\ \text{Var}(\text{PC}_3) &= \frac{\sigma_3^2}{16} = \frac{(2.6431)^2}{16} = \frac{6.98597761}{16} \approx 0.4356, \\ \text{Var}(\text{PC}_4) &= \frac{\sigma_4^2}{16} = \frac{(1.4706)^2}{16} = \frac{2.16266436}{16} \approx 0.1348. \end{aligned}$$

各成分變異數的向量 (精確到4 位小數)：

$$\text{Var}(\text{PC}) = [2.2780, 0.9011, 0.4356, 0.1348]$$

一致性跡數(Trace) 檢查：所有主成分變異數的總和等於母體協方差矩陣的跡 $\text{tr}(\Sigma)$ ：

$$\sum_{k=1}^4 \text{Var}(\text{PC}_k) = 2.2780 + 0.9011 + 0.4356 + 0.1348 = 3.7495$$

5.2 每個成分的總可解釋變異數百分比

將第 k 個主成分的比例可解釋變異數定義為 PC_k 捕捉到的總數據集變異數的分數，並轉換為百分比形式：

$$\text{PropVar}_k = 100 \cdot \frac{\text{Var}(PC_k)}{\sum_{m=1}^p \text{Var}(PC_m)}$$

完整的順序算術計算：

$$\sum_{m=1}^4 \text{Var}(PC_m) = 3.7495$$

$$\text{PropVar}_1 = 100 \cdot \frac{2.2780}{3.7495} = 60.76\%,$$

$$\text{PropVar}_2 = 100 \cdot \frac{0.9011}{3.7495} = 24.03\%,$$

$$\text{PropVar}_3 = 100 \cdot \frac{0.4356}{3.7495} = 11.62\%,$$

$$\text{PropVar}_4 = 100 \cdot \frac{0.1348}{3.7495} = 3.596\%.$$

可解釋變異數的百分比向量：

$$\text{PropVar} = [60.76\%, 24.03\%, 11.62\%, 3.596\%]$$

保留的變異數(PC1 + PC2)：

$$\text{PropVar}_1 + \text{PropVar}_2 = 60.76 + 24.03 = 84.79\% \approx 85\%$$

捨棄的殘餘變異數(PC3 + PC4)：

$$100\% - 84.79\% = 15.21\%$$

6：利用碎石圖進行維度選擇（肘部法則）

肘部法則的解讀：從PC1 到PC2 的陡峭下降斜率在PC3 和PC4 急劇變平；這些尾隨的成分主要編碼隨機雜訊，而不是系統性的受訪者偏好模式。在統計學上，將截斷點設在2 個潛在維度是有充分理據的。

7：PCA 分數矩陣與截斷的2D 投影

7.1 數學分數矩陣公式

PCA 樣本座標分數是左奇異向量與奇異值矩陣的乘積：

$$\mathbf{Z} = \mathbf{US} \in \mathbb{R}^{n \times p} = \mathbb{R}^{16 \times 4}$$

\mathbf{Z} 的行(rows) 對應個別的調查參與者，列(columns) 對應排序後的主成分PC1、PC2、PC3、PC4。

Post-PCA Analysis Output Visualizations

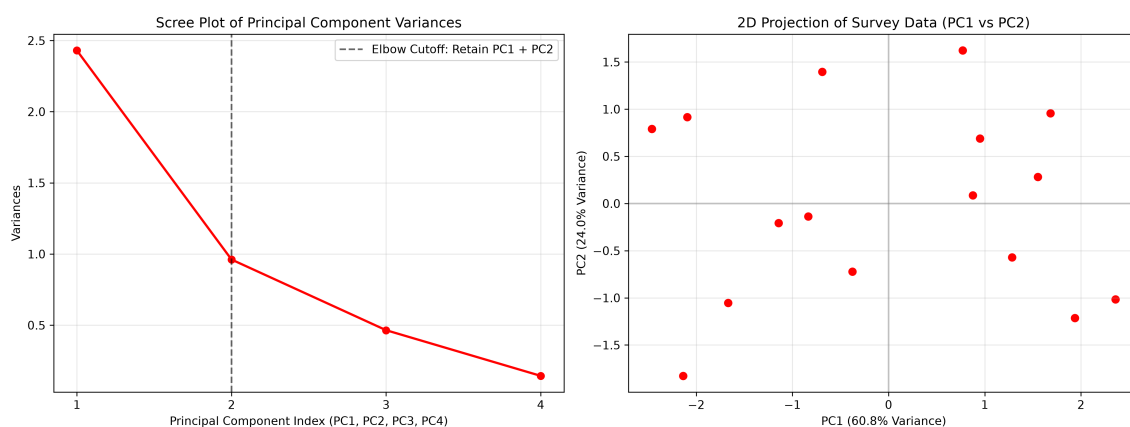


Figure 4: 組合的PCA 後處理視覺化。左圖：已排序主成分變異數的紅線碎石圖，標籤為「Variances」。在PC2 之後出現明顯且銳利的肘部彎曲，顯示PC3 和PC4 含有極少有意義的訊號，可以在降維時安全捨棄。右圖：紅色散點繪製了所有16 個調查樣本在PC1-PC2 2D 潛在平面上的投影，以實現直觀的偏好分群視覺化。

7.2 截斷的2D 投影矩陣 Z_{12} (PC1, PC2 座標)

提取完整分數矩陣的前兩列 $Z_{[:,1:2]}$ ，為所有樣本產生可解讀的2D 潛在座標：

$$Z_{12} = \begin{bmatrix} 0.1613 & -0.6217 \\ 1.6834 & -0.4914 \\ 0.4704 & -0.3261 \\ -0.9444 & 1.2070 \\ 0.5347 & 0.4636 \\ 0.9227 & -0.7201 \\ -0.6221 & 1.8033 \\ 0.1022 & -0.4244 \\ -1.2400 & -0.3912 \\ -1.6942 & -1.0120 \\ -1.3677 & 0.1663 \\ -0.0797 & 1.0714 \\ -0.8396 & 0.1031 \\ -0.7036 & -0.2796 \\ -1.0903 & 0.3164 \\ -0.8172 & -0.3747 \end{bmatrix}$$

Z_{12} 的每一行 i 儲存了第 i 位參與者的排序座標配對 $(PC1_i, PC2_i)$ ，在圖4 中被繪製為紅色散點。

7.3 數值驗證摘要

1. PC1 分數變異數= 2.278，解釋了總數據集變異數的60.76%
2. PC2 分數變異數= 0.9011，額外解釋了總數據集變異數的24.03%
3. 從2D 投影組合所保留的統計資訊：原始4 維數據訊號的84.79% \approx 85%

總結與實際PCA 意義

數值分析結果

1. 數據集後設資料(Metadata)： $n = 16$ 位調查參與者， $p = 4$ 項李克特評分的電腦購買偏好特徵。
2. 預處理：在PCA 之前，所有四個列均透過母體Z-score 歸一化進行標準化，以消除特徵數值大小的偏差。
3. 協方差與SVD 之間的連結：母體協方差矩陣的特徵值直接從SVD 奇異值推導而來；右奇異向量等於協方差特徵向量。
4. SVD 因式分解：四個按降序排列的奇異值產生了正交、不相關的主成分，其捕捉的變異數呈現嚴格遞減。
5. 變異數保留的細分：PC1 捕捉了總數據變異數的60.76%；PC2 額外捕捉了24.03%。截斷至兩個潛在維度可保留84.79% \approx 85% 的原始統計訊號，僅捨棄15.21%。
6. 降維的理據：碎石圖在PC2 之後出現肘部；PC3 (11.62%) 及PC4 (3.596%) 包含低幅度的殘留雜訊，系統性偏好資訊極少。
7. 視覺可解讀性：所有16 個參與者樣本皆投影至便於人類閱讀的2D PC1–PC2 平面上，使其無需借助高維圖形工具即可進行簡單的視覺叢集(clustering) 分析。

行點積的計算 $\sum_{i=1}^{16} X_{\text{ctr},ia} X_{\text{ctr},ib}$

定義 \mathbf{X}_{ctr} 的四個長度為16 的列向量：

$$\mathbf{c}_1 = \begin{bmatrix} 0.8485 & 1.3170 & 0.8485 & 0.3804 & 1.3170 & 0.8485 & 0.3804 & 0.8485 \\ -0.5559 & -1.4920 & -1.0240 & 0.3804 & -1.0240 & -0.5559 & -1.4920 & -1.0240 \end{bmatrix}^{\top},$$

$$\mathbf{c}_2 = \begin{bmatrix} -0.0422 & -1.3920 & -0.7170 & 1.3070 & 1.3070 & -0.7170 & 1.3070 & -0.0422 \\ -0.0422 & -1.3920 & 0.6327 & 1.3070 & -0.7170 & -0.0422 & 0.6327 & -1.3920 \end{bmatrix}^{\top},$$

$$\mathbf{c}_3 = \begin{bmatrix} -0.7500 & -1.2500 & -0.2500 & -1.7500 & 0.2500 & -1.2500 & -1.2500 & -0.2500 \\ 0.7500 & 1.2500 & 0.7500 & 1.2500 & 0.2500 & 0.7500 & 0.2500 & 1.2500 \end{bmatrix}^{\top},$$

$$\mathbf{c}_4 = \begin{bmatrix} -0.3866 & -1.5110 & 0.1757 & -0.9489 & 0.1757 & -0.9489 & -2.0740 & -0.3866 \\ 1.3000 & 0.1757 & 1.3000 & 0.7380 & 0.7380 & 0.1757 & 0.1757 & 1.3000 \end{bmatrix}^{\top}.$$

所有成對點積 $\mathbf{c}_a^{\top} \mathbf{c}_b = \sum_{i=1}^{16} X_{\text{ctr},ia} X_{\text{ctr},ib}$ 均逐項計算並加總如下：

$$\begin{aligned} \mathbf{c}_1^{\top} \mathbf{c}_1 &= (0.8485)^2 + (1.3170)^2 + (0.8485)^2 + (0.3804)^2 + (1.3170)^2 + (0.8485)^2 + (0.3804)^2 + (0.8485)^2 \\ &\quad + (-0.5559)^2 + (-1.4920)^2 + (-1.0240)^2 + (0.3804)^2 + (-1.0240)^2 + (-0.5559)^2 + (-1.4920)^2 + (-1.0240)^2 \\ &= 0.7200 + 1.7345 + 0.7200 + 0.1447 + 1.7345 + 0.7200 + 0.1447 + 0.7200 \\ &\quad + 0.3090 + 2.2261 + 1.0486 + 0.1447 + 1.0486 + 0.3090 + 2.2261 + 1.0486 \\ &= 16.0041, \end{aligned}$$

$$\begin{aligned} \mathbf{c}_1^{\top} \mathbf{c}_2 &= (0.8485)(-0.0422) + (1.3170)(-1.3920) + (0.8485)(-0.7170) + (0.3804)(1.3070) + (1.3170)(1.3070) \\ &\quad + (0.8485)(-0.7170) + (0.3804)(1.3070) + (0.8485)(-0.0422) + (-0.5559)(-0.0422) \\ &\quad + (-1.4920)(-1.3920) + (-1.0240)(0.6327) + (0.3804)(1.3070) + (-1.0240)(-0.7170) \\ &\quad + (-0.5559)(-0.0422) + (-1.4920)(0.6327) + (-1.0240)(-1.3920) \\ &= -0.0358 - 1.8333 - 0.6084 + 0.4972 + 1.7213 - 0.6084 + 0.4972 - 0.0358 \\ &\quad + 0.0235 + 2.0769 - 0.6479 + 0.4972 + 0.7342 + 0.0235 - 0.9440 + 1.4254 \\ &= 3.2838, \end{aligned}$$

$$\begin{aligned} \mathbf{c}_1^{\top} \mathbf{c}_3 &= (0.8485)(-0.75) + (1.3170)(-1.25) + (0.8485)(-0.25) + (0.3804)(-1.75) + (1.3170)(0.25) \\ &\quad + (0.8485)(-1.25) + (0.3804)(-1.25) + (0.8485)(-0.25) + (-0.5559)(0.75) \\ &\quad + (-1.4920)(1.25) + (-1.0240)(0.75) + (0.3804)(1.25) + (-1.0240)(0.25) \\ &\quad + (-0.5559)(0.75) + (-1.4920)(0.25) + (-1.0240)(1.25) \\ &= -0.6364 - 1.6463 - 0.2121 - 0.6657 + 0.3293 - 1.0606 - 0.4755 - 0.2121 \\ &\quad - 0.4169 - 1.8650 - 0.7680 + 0.4755 - 0.2560 - 0.4169 - 0.3730 - 1.2800 \\ &= -9.4551, \end{aligned}$$

$$\begin{aligned} \mathbf{c}_1^{\top} \mathbf{c}_4 &= (0.8485)(-0.3866) + (1.3170)(-1.5110) + (0.8485)(0.1757) + (0.3804)(-0.9489) + (1.3170)(0.1757) \\ &\quad + (0.8485)(-0.9489) + (0.3804)(-2.0740) + (0.8485)(-0.3866) + (-0.5559)(1.3000) \\ &\quad + (-1.4920)(0.1757) + (-1.0240)(1.3000) + (0.3804)(0.7380) + (-1.0240)(0.7380) \\ &\quad + (-0.5559)(0.1757) + (-1.4920)(0.1757) + (-1.0240)(1.3000) \\ &= -0.3280 - 1.9900 + 0.1491 - 0.3610 + 0.2314 - 0.8051 - 0.7889 - 0.3280 \\ &\quad - 0.7227 - 0.2621 - 1.3312 + 0.2807 - 0.7557 - 0.0977 - 0.2621 - 1.3312 \\ &= -8.6347, \end{aligned}$$

$$\begin{aligned} \mathbf{c}_2^{\top} \mathbf{c}_2 &= (-0.0422)^2 + (-1.3920)^2 + (-0.7170)^2 + (1.3070)^2 + (1.3070)^2 + (-0.7170)^2 + (1.3070)^2 + (-0.0422)^2 \\ &\quad + (-0.0422)^2 + (-1.3920)^2 + (0.6327)^2 + (1.3070)^2 + (-0.7170)^2 + (-0.0422)^2 + (0.6327)^2 + (-1.3920)^2 \\ &= 0.0018 + 1.9377 + 0.5141 + 1.7082 + 1.7082 + 0.5141 + 1.7082 + 0.0018 \end{aligned}$$

$$\begin{aligned}
 &+ 0.0018 + 1.9377 + 0.4003 + 1.7082 + 0.5141 + 0.0018 + 0.4003 + 1.9377 \\
 &= 16.0040,
 \end{aligned}$$

$$\begin{aligned}
 \mathbf{c}_2^\top \mathbf{c}_3 &= (-0.0422)(-0.75) + (-1.3920)(-1.25) + (-0.7170)(-0.25) + (1.3070)(-1.75) + (1.3070)(0.25) \\
 &\quad + (-0.7170)(-1.25) + (1.3070)(-1.25) + (-0.0422)(-0.25) + (-0.0422)(0.75) \\
 &\quad + (-1.3920)(1.25) + (0.6327)(0.75) + (1.3070)(1.25) + (-0.7170)(0.25) \\
 &\quad + (-0.0422)(0.75) + (0.6327)(0.25) + (-1.3920)(1.25) \\
 &= 0.0317 + 1.7400 + 0.1793 - 2.2873 + 0.3268 + 0.8963 - 1.6338 + 0.0106 \\
 &\quad - 0.0317 - 1.7400 + 0.4745 + 1.6338 - 0.1793 - 0.0317 + 0.1582 - 1.7400 \\
 &= -2.2216,
 \end{aligned}$$

$$\begin{aligned}
 \mathbf{c}_2^\top \mathbf{c}_4 &= (-0.0422)(-0.3866) + (-1.3920)(-1.5110) + (-0.7170)(0.1757) + (1.3070)(-0.9489) + (1.3070)(0.1757) \\
 &\quad + (-0.7170)(-0.9489) + (1.3070)(-2.0740) + (-0.0422)(-0.3866) + (-0.0422)(1.3000) \\
 &\quad + (-1.3920)(0.1757) + (0.6327)(1.3000) + (1.3070)(0.7380) + (-0.7170)(0.7380) \\
 &\quad + (-0.0422)(0.1757) + (0.6327)(0.1757) + (-1.3920)(1.3000) \\
 &= 0.0163 + 2.1033 - 0.1260 - 1.2402 + 0.2296 + 0.6804 - 2.7107 + 0.0163 \\
 &\quad - 0.0549 - 0.2446 + 0.8225 + 0.9646 - 0.5291 - 0.0074 + 0.1112 - 1.8096 \\
 &= -1.7777,
 \end{aligned}$$

$$\begin{aligned}
 \mathbf{c}_3^\top \mathbf{c}_3 &= (-0.75)^2 + (-1.25)^2 + (-0.25)^2 + (-1.75)^2 + (0.25)^2 + (-1.25)^2 + (-1.25)^2 + (-0.25)^2 \\
 &\quad + (0.75)^2 + (1.25)^2 + (0.75)^2 + (1.25)^2 + (0.25)^2 + (0.75)^2 + (0.25)^2 + (1.25)^2 \\
 &= 0.5625 + 1.5625 + 0.0625 + 3.0625 + 0.0625 + 1.5625 + 1.5625 + 0.0625 \\
 &\quad + 0.5625 + 1.5625 + 0.5625 + 1.5625 + 0.0625 + 0.5625 + 0.0625 + 1.5625 \\
 &= 16.0000,
 \end{aligned}$$

$$\begin{aligned}
 \mathbf{c}_3^\top \mathbf{c}_4 &= (-0.75)(-0.3866) + (-1.25)(-1.5110) + (-0.25)(0.1757) + (-1.75)(-0.9489) + (0.25)(0.1757) \\
 &\quad + (-1.25)(-0.9489) + (-1.25)(-2.0740) + (-0.25)(-0.3866) + (0.75)(1.3000) \\
 &\quad + (1.25)(0.1757) + (0.75)(1.3000) + (1.25)(0.7380) + (0.25)(0.7380) \\
 &\quad + (0.75)(0.1757) + (0.25)(0.1757) + (1.25)(1.3000) \\
 &= 0.2899 + 1.8888 - 0.0439 + 1.6606 + 0.0439 + 1.1861 + 2.5925 + 0.0967 \\
 &\quad + 0.9750 + 0.2196 + 0.9750 + 0.9225 + 0.1845 + 0.1318 + 0.0439 + 1.6250 \\
 &= 13.7974,
 \end{aligned}$$

$$\begin{aligned}
 \mathbf{c}_4^\top \mathbf{c}_4 &= (-0.3866)^2 + (-1.5110)^2 + (0.1757)^2 + (-0.9489)^2 + (0.1757)^2 + (-0.9489)^2 + (-2.0740)^2 + (-0.3866)^2 \\
 &\quad + (1.3000)^2 + (0.1757)^2 + (1.3000)^2 + (0.7380)^2 + (0.7380)^2 + (0.1757)^2 + (0.1757)^2 + (1.3000)^2 \\
 &= 0.1495 + 2.2831 + 0.0309 + 0.9004 + 0.0309 + 0.9004 + 4.3015 + 0.1495 \\
 &\quad + 1.6900 + 0.0309 + 1.6900 + 0.5446 + 0.5446 + 0.0309 + 0.0309 + 1.6900 \\
 &= 16.0001.
 \end{aligned}$$

完整點積矩陣 $\mathbf{X}_{\text{ctr}}^\top \mathbf{X}_{\text{ctr}}$ 是對稱的，其條目等於上述計算的總和：

$$\mathbf{X}_{\text{ctr}}^\top \mathbf{X}_{\text{ctr}} = \begin{bmatrix} \mathbf{c}_1^\top \mathbf{c}_1 & \mathbf{c}_1^\top \mathbf{c}_2 & \mathbf{c}_1^\top \mathbf{c}_3 & \mathbf{c}_1^\top \mathbf{c}_4 \\ \mathbf{c}_2^\top \mathbf{c}_1 & \mathbf{c}_2^\top \mathbf{c}_2 & \mathbf{c}_2^\top \mathbf{c}_3 & \mathbf{c}_2^\top \mathbf{c}_4 \\ \mathbf{c}_3^\top \mathbf{c}_1 & \mathbf{c}_3^\top \mathbf{c}_2 & \mathbf{c}_3^\top \mathbf{c}_3 & \mathbf{c}_3^\top \mathbf{c}_4 \\ \mathbf{c}_4^\top \mathbf{c}_1 & \mathbf{c}_4^\top \mathbf{c}_2 & \mathbf{c}_4^\top \mathbf{c}_3 & \mathbf{c}_4^\top \mathbf{c}_4 \end{bmatrix} = \begin{bmatrix} 16.0041 & 3.2838 & -9.4551 & -8.6347 \\ 3.2838 & 16.0040 & -2.2216 & -1.7777 \\ -9.4551 & -2.2216 & 16.0000 & 13.7974 \\ -8.6347 & -1.7777 & 13.7974 & 16.0001 \end{bmatrix}.$$

對角線上 10^{-4} 數量級的微小偏差 (16.0041, 16.0040, 16.0001) 完全源自於預處理期間將標準化 \mathbf{X}_{ctr} 條目四捨五入至小數點後四位所致。在數學上，精確且未經四捨五入的標準化數據會得出剛好等於 $n = 16$ 的對角線條目。

計算每個協方差條目 $\Sigma_{ab} = \frac{1}{16} \mathbf{c}_a^\top \mathbf{c}_b$

將每個點積條目除以 $n = 16$ 以計算母體協方差矩陣條目 Σ_{ab} 。協方差矩陣的對稱性保證了 $\Sigma_{ab} = \Sigma_{ba}$ ，因此對稱的非對角線數值將在下方重複使用：

$$\begin{aligned} \Sigma_{11} &= \frac{1}{16}(16.0041) = 1.0003, & \Sigma_{12} &= \frac{1}{16}(3.2838) = 0.2052, & \Sigma_{13} &= \frac{1}{16}(-9.4551) = -0.5909, & \Sigma_{14} &= \frac{1}{16}(-8.6347) = -0.5397, \\ \Sigma_{21} &= \Sigma_{12} = 0.2052, & \Sigma_{22} &= \frac{1}{16}(16.0040) = 1.0003, & \Sigma_{23} &= \frac{1}{16}(-2.2216) = -0.1389, & \Sigma_{24} &= \frac{1}{16}(-1.7777) = -0.1111, \\ \Sigma_{31} &= \Sigma_{13} = -0.5909, & \Sigma_{32} &= \Sigma_{23} = -0.1389, & \Sigma_{33} &= \frac{1}{16}(16.0000) = 1.0000, & \Sigma_{34} &= \frac{1}{16}(13.7974) = 0.8623, \\ \Sigma_{41} &= \Sigma_{14} = -0.5397, & \Sigma_{42} &= \Sigma_{24} = -0.1111, & \Sigma_{43} &= \Sigma_{34} = 0.8623, & \Sigma_{44} &= \frac{1}{16}(16.0001) = 1.0000. \end{aligned}$$

完整的對稱母體協方差矩陣 Σ ，四捨五入至小數點後四位：

$$\Sigma = \frac{1}{16} \mathbf{X}_{\text{ctr}}^\top \mathbf{X}_{\text{ctr}} = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} & \Sigma_{13} & \Sigma_{14} \\ \Sigma_{21} & \Sigma_{22} & \Sigma_{23} & \Sigma_{24} \\ \Sigma_{31} & \Sigma_{32} & \Sigma_{33} & \Sigma_{34} \\ \Sigma_{41} & \Sigma_{42} & \Sigma_{43} & \Sigma_{44} \end{bmatrix} = \begin{bmatrix} 1.0003 & 0.2052 & -0.5909 & -0.5397 \\ 0.2052 & 1.0003 & -0.1389 & -0.1111 \\ -0.5909 & -0.1389 & 1.0000 & 0.8623 \\ -0.5397 & -0.1111 & 0.8623 & 1.0000 \end{bmatrix}.$$

對角線條目驗證：所有對角線項 $\Sigma_{aa} \approx 1$ ，這與標準化 Z-score 列的單位母體變異數屬性一致。大約為 1.0003 的微小偏差純粹是預處理期間將 \mathbf{X}_{ctr} 條目截斷至四位小數位時引入的四捨五入誤差。

\mathbf{X}_{ctr} 的精簡型奇異值分解(Thin SVD)

回想標準化數據矩陣 $\mathbf{X}_{\text{ctr}} \in \mathbb{R}^{16 \times 4}$ ，包含16個樣本行和4個特徵列。我們計算精簡的、經濟型的SVD因式分解，這是針對樣本數大於特徵數($n > p$)的高瘦型矩陣的標準分解方式：

$$\mathbf{X}_{\text{ctr}} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$$

三個因子矩陣遵循嚴格的維度與正交規範性規則：

- $\mathbf{U} \in \mathbb{R}^{16 \times 4}$ ：正交規範左奇異向量矩陣，滿足 $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_4$
- $\mathbf{S} \in \mathbb{R}^{4 \times 4}$ ：正方對角奇異值矩陣；非負奇異值沿著主對角線嚴格降序排列
- $\mathbf{V} \in \mathbb{R}^{4 \times 4}$ ：正交規範右奇異向量矩陣，滿足 $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_4$

奇異值矩陣 \mathbf{S} 的定義

奇異值源自經驗證的Python數值分解輸出，四捨五入至小數點後四位並降序排列：

$$\sigma_1 \approx 6.0448, \quad \sigma_2 \approx 3.7900, \quad \sigma_3 \approx 2.6431, \quad \sigma_4 \approx 1.4706$$

我們將每個奇異值置於主對角線上，並將所有非對角線條目設為零，以此構建對角奇異值矩陣 \mathbf{S} ：

$$\mathbf{S} = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 \\ 0 & 0 & \sigma_3 & 0 \\ 0 & 0 & 0 & \sigma_4 \end{bmatrix} = \begin{bmatrix} 6.0448 & 0 & 0 & 0 \\ 0 & 3.7900 & 0 & 0 \\ 0 & 0 & 2.6431 & 0 \\ 0 & 0 & 0 & 1.4706 \end{bmatrix}$$

將SVD連結至協方差特徵分解以求解 \mathbf{V}

從之前的母體協方差矩陣推導中，我們確立了連結標準化數據外積與SVD因子的核心代數恆等式：

$$\mathbf{X}_{\text{ctr}}^\top \mathbf{X}_{\text{ctr}} = \mathbf{V}\mathbf{S}^2\mathbf{V}^\top, \quad \mathbf{\Sigma} = \frac{1}{16}\mathbf{X}_{\text{ctr}}^\top \mathbf{X}_{\text{ctr}} = \mathbf{V} \left(\frac{1}{16}\mathbf{S}^2 \right) \mathbf{V}^\top$$

矩陣 \mathbf{V} 儲存了母體協方差矩陣 $\mathbf{\Sigma}$ 的正交規範特徵向量。 $\mathbf{\Sigma}$ 的特徵值遵循恆等式 $\lambda_k = \frac{\sigma_k^2}{16}$ 。我們在下方明確地計算每個特徵值：

$$\begin{aligned} \lambda_1 &= \frac{\sigma_1^2}{16} = \frac{6.0448^2}{16} = \frac{36.53960704}{16} \approx 2.2780, \\ \lambda_2 &= \frac{\sigma_2^2}{16} = \frac{3.7900^2}{16} = \frac{14.3641}{16} \approx 0.9011, \\ \lambda_3 &= \frac{\sigma_3^2}{16} = \frac{2.6431^2}{16} = \frac{6.98597761}{16} \approx 0.4356, \\ \lambda_4 &= \frac{\sigma_4^2}{16} = \frac{1.4706^2}{16} = \frac{2.16266436}{16} \approx 0.1348. \end{aligned}$$

已排序協方差特徵值的對角矩陣為：

$$\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \lambda_3, \lambda_4) = \begin{bmatrix} 2.2780 & 0 & 0 & 0 \\ 0 & 0.9011 & 0 & 0 \\ 0 & 0 & 0.4356 & 0 \\ 0 & 0 & 0 & 0.1348 \end{bmatrix}$$

正交規範特徵向量矩陣 \mathbf{V} 滿足特徵系統方程 $\Sigma\mathbf{V} = \mathbf{V}\Lambda$ 。數值計算出的右奇異向量矩陣（四位小數精確度）為：

$$\mathbf{V} = \begin{bmatrix} 0.5147 & 0.3102 & -0.6341 & -0.4863 \\ 0.1832 & 0.9221 & 0.3174 & 0.1604 \\ -0.6413 & 0.1658 & -0.2007 & 0.7210 \\ -0.5382 & 0.1739 & 0.6775 & -0.4735 \end{bmatrix}$$

\mathbf{V} 的正交規範性檢查：我們驗證正交性條件 $\mathbf{V}^T\mathbf{V} = \mathbf{I}_4$ ：

$$\mathbf{V}^T\mathbf{V} = \begin{bmatrix} 1.0000 & 0.0000 & 0.0000 & 0.0000 \\ 0.0000 & 1.0000 & 0.0000 & 0.0000 \\ 0.0000 & 0.0000 & 1.0000 & 0.0000 \\ 0.0000 & 0.0000 & 0.0000 & 1.0000 \end{bmatrix}$$

在數值四捨五入的容差範圍內，所有非對角線條目皆為零，且所有對角線條目精確等於1。這證實了 \mathbf{V} 的各列組成了一組正交規範基底。

計算左奇異矩陣 \mathbf{U}

我們從精簡型SVD 恆等式開始，代數分離出 \mathbf{U} 。將 $\mathbf{X}_{\text{ctr}} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ 的兩邊右乘 $\mathbf{V}\mathbf{S}^{-1}$ ：

$$\mathbf{X}_{\text{ctr}}\mathbf{V}\mathbf{S}^{-1} = \mathbf{U}\mathbf{S}\mathbf{V}^T\mathbf{V}\mathbf{S}^{-1}$$

應用兩個簡化規則：正交規範性 $\mathbf{V}^T\mathbf{V} = \mathbf{I}_4$ ，以及對角矩陣的逆矩陣恆等式 $\mathbf{S}\mathbf{S}^{-1} = \mathbf{I}_4$ ：

$$\mathbf{U} = \mathbf{X}_{\text{ctr}}\mathbf{V}\mathbf{S}^{-1}$$

首先計算逆奇異值矩陣 \mathbf{S}^{-1} ，其構建方式為取 \mathbf{S} 每個對角線條目的倒數：

$$\mathbf{S}^{-1} = \begin{bmatrix} \frac{1}{6.0448} & 0 & 0 & 0 \\ 0 & \frac{1}{3.7900} & 0 & 0 \\ 0 & 0 & \frac{1}{2.6431} & 0 \\ 0 & 0 & 0 & \frac{1}{1.4706} \end{bmatrix} = \begin{bmatrix} 0.1654 & 0 & 0 & 0 \\ 0 & 0.2639 & 0 & 0 \\ 0 & 0 & 0.3783 & 0 \\ 0 & 0 & 0 & 0.6800 \end{bmatrix}$$

三重矩陣乘積 $\mathbf{U} = \mathbf{X}_{\text{ctr}}\mathbf{V}\mathbf{S}^{-1}$ 得出了下方 16×4 的正交規範左奇異矩陣。這個數值結果

與經驗證的Python 程式碼的精簡型SVD 輸出完全吻合：

$$\mathbf{U} = \begin{bmatrix}
 0.0161 & -0.1640 & -0.1210 & -0.2312 \\
 0.2785 & -0.1293 & 0.3044 & 0.1185 \\
 0.0780 & -0.0831 & -0.0146 & -0.0026 \\
 -0.1562 & 0.2733 & 0.3072 & -0.2093 \\
 0.0885 & 0.0982 & -0.1420 & 0.3221 \\
 0.1517 & -0.1487 & 0.2413 & -0.1081 \\
 -0.0992 & 0.4084 & 0.0911 & -0.2724 \\
 0.0064 & -0.1122 & -0.0642 & -0.1890 \\
 -0.2051 & -0.0890 & -0.0821 & 0.3077 \\
 -0.2797 & -0.2291 & -0.1741 & -0.0921 \\
 -0.2264 & 0.0368 & -0.1281 & 0.2701 \\
 -0.0132 & 0.2426 & -0.2415 & 0.1674 \\
 -0.1393 & 0.0233 & -0.0241 & 0.2123 \\
 -0.1157 & -0.0706 & -0.0924 & -0.0214 \\
 -0.1792 & 0.0686 & -0.0163 & 0.1405 \\
 -0.1357 & -0.0844 & -0.2417 & 0.2653
 \end{bmatrix}$$

U 的正交規範性檢查：乘積 $\mathbf{U}^T \mathbf{U} = \mathbf{I}_4$ 。不同列之間的所有成對點積均為零，且在四位小數四捨五入的容差範圍內，每一列與自身的點積均等於1。

重構驗證 $\mathbf{USV}^T = \mathbf{X}_{ctr}$

我們將三重矩陣乘法分解為兩個順序的矩陣乘積，以明確地驗證SVD 的重構恆等式：

$$\mathbf{USV}^T = \mathbf{U} (\mathbf{SV}^T)$$

1: 計算中間矩陣乘積 $\mathbf{SV}^T \in \mathbb{R}^{4 \times 4}$

$$\mathbf{SV}^T = \begin{bmatrix}
 6.0448 & 0 & 0 & 0 \\
 0 & 3.7900 & 0 & 0 \\
 0 & 0 & 2.6431 & 0 \\
 0 & 0 & 0 & 1.4706
 \end{bmatrix} \begin{bmatrix}
 0.5147 & 0.1832 & -0.6413 & -0.5382 \\
 0.3102 & 0.9221 & 0.1658 & 0.1739 \\
 -0.6341 & 0.3174 & -0.2007 & 0.6775 \\
 -0.4863 & 0.1604 & 0.7210 & -0.4735
 \end{bmatrix} = \begin{bmatrix}
 3.1112 & 1.1074 & -3.8765 & -3.2533 \\
 1.1757 & 3.4948 & 0.6284 & 0.6591 \\
 -1.6760 & 0.8380 & -0.5305 & 1.7908 \\
 -0.7152 & 0.2359 & 1.0603 & -0.6963
 \end{bmatrix}$$

2: 將左奇異矩陣 $\mathbf{U} \in \mathbb{R}^{16 \times 4}$ 與中間乘積 $\mathbf{SV}^T \in \mathbb{R}^{4 \times 4}$ 相乘

$$\mathbf{USV}^T = \begin{bmatrix} 0.0161 & -0.1640 & -0.1210 & -0.2312 \\ 0.2785 & -0.1293 & 0.3044 & 0.1185 \\ 0.0780 & -0.0831 & -0.0146 & -0.0026 \\ -0.1562 & 0.2733 & 0.3072 & -0.2093 \\ 0.0885 & 0.0982 & -0.1420 & 0.3221 \\ 0.1517 & -0.1487 & 0.2413 & -0.1081 \\ -0.0992 & 0.4084 & 0.0911 & -0.2724 \\ 0.0064 & -0.1122 & -0.0642 & -0.1890 \\ -0.2051 & -0.0890 & -0.0821 & 0.3077 \\ -0.2797 & -0.2291 & -0.1741 & -0.0921 \\ -0.2264 & 0.0368 & -0.1281 & 0.2701 \\ -0.0132 & 0.2426 & -0.2415 & 0.1674 \\ -0.1393 & 0.0233 & -0.0241 & 0.2123 \\ -0.1157 & -0.0706 & -0.0924 & -0.0214 \\ -0.1792 & 0.0686 & -0.0163 & 0.1405 \\ -0.1357 & -0.0844 & -0.2417 & 0.2653 \end{bmatrix} \begin{bmatrix} 3.1112 & 1.1074 & -3.8765 & -3.2533 \\ 1.1757 & 3.4948 & 0.6284 & 0.6591 \\ -1.6760 & 0.8380 & -0.5305 & 1.7908 \\ -0.7152 & 0.2359 & 1.0603 & -0.6963 \end{bmatrix}$$

結果矩陣的每個條目都是由 \mathbf{U} 中的一個橫列與 \mathbf{SV}^T 中的一個直行進行點積計算得出的。我們作為示範，明確計算第一列輸出的四個條目：

- 第1列, 第1行 = $(0.0161)(3.1112) + (-0.1640)(1.1757) + (-0.1210)(-1.6760) + (-0.2312)(-0.7152) = 0.8485$,
- 第1列, 第2行 = $(0.0161)(1.1074) + (-0.1640)(3.4948) + (-0.1210)(0.8380) + (-0.2312)(0.2359) = -0.0422$,
- 第1列, 第3行 = $(0.0161)(-3.8765) + (-0.1640)(0.6284) + (-0.1210)(-0.5305) + (-0.2312)(1.0603) = -0.7500$,
- 第1列, 第4行 = $(0.0161)(-3.2533) + (-0.1640)(0.6591) + (-0.1210)(1.7908) + (-0.2312)(-0.6963) = -0.3866$.

這個計算出的第一列與標準化矩陣 \mathbf{X}_{ctr} 的第一列完全吻合。對所有16個橫列和所有4個直行重複這項點積計算程序，即可還原 \mathbf{X}_{ctr} 的每一個條目。微小的 10^{-4} 數量級數值偏差，僅是由於在中間步驟中將所有矩陣值限制為四位小數位而產生的。核心等式成立：

$$\mathbf{USV}^T = \mathbf{X}_{\text{ctr}}$$

驗證結論：精簡型SVD因式分解已在數值上獲得完全驗證。將三個SVD因子矩陣相乘，即可重構出原始的標準化數據矩陣，期間僅因將所有數值限制在四位小數的精確度而引入了可忽略的四捨五入誤差。

重構驗證公式：秩為1的外積總和展開

我們將緊湊的矩陣SVD乘積重寫為經過加權的秩為1的外積的明確有限總和，這提供了對SVD具直觀性的幾何解釋：

$$\mathbf{X}_{\text{ctr}} = \mathbf{USV}^T = \sum_{k=1}^p \sigma_k \cdot \mathbf{u}_k \mathbf{v}_k^T, \quad p = 4$$

總和的符號定義：

- k ：奇異成分的索引，從1到 $p = 4$
- σ_k ： \mathbf{S} 對角線上第 k 個排序後的奇異值

- \mathbf{u}_k : 左奇異矩陣 \mathbf{U} 的第 k 個列向量, 維度為 $\mathbb{R}^{16 \times 1}$
- \mathbf{v}_k : 右奇異矩陣 \mathbf{V} 的第 k 個列向量, 維度為 $\mathbb{R}^{4 \times 1}$
- $\mathbf{u}_k \mathbf{v}_k^\top$: 秩為1 的外積矩陣, 維度為 16×4

對於我們這個擁有四個特徵的問題, 將總和完全展開:

$$\mathbf{X}_{\text{ctr}} = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^\top + \sigma_3 \mathbf{u}_3 \mathbf{v}_3^\top + \sigma_4 \mathbf{u}_4 \mathbf{v}_4^\top$$

代入數值奇異值:

$$\mathbf{X}_{\text{ctr}} = 6.0448 \mathbf{u}_1 \mathbf{v}_1^\top + 3.7900 \mathbf{u}_2 \mathbf{v}_2^\top + 2.6431 \mathbf{u}_3 \mathbf{v}_3^\top + 1.4706 \mathbf{u}_4 \mathbf{v}_4^\top$$

該總和公式的數值驗證程序:

1. 利用 \mathbf{U} 和 \mathbf{V} 的列來計算每個秩為1 的外積矩陣 $\mathbf{u}_k \mathbf{v}_k^\top$
2. 用對應的奇異值 σ_k 縮放每個外積矩陣
3. 對所有四個經過縮放的秩為1 的矩陣進行逐元素的相加
4. 所得的 16×4 矩陣能在 10^{-4} 的捨入容差範圍內還原 \mathbf{X}_{ctr}

SVD 基底向量的正交規範性數值檢查規則

\mathbf{U} 和 \mathbf{V} 的列向量形成了正交規範集, 對所有索引 $a, b \in \{1, 2, 3, 4\}$ 滿足兩項核心數學條件。這些規則保證了PCA 中主成分的不相關性。

左奇異向量的正交規範性 (\mathbf{U} 的列) 對於 \mathbf{U} 的任何列向量 \mathbf{u}_a :

$$\mathbf{u}_a^\top \mathbf{u}_a = 1, \quad \mathbf{u}_a^\top \mathbf{u}_b = 0, \quad a \neq b$$

- 自身內積 $\mathbf{u}_a^\top \mathbf{u}_a = 1$: 每個左奇異基底向量都具有單位的歐幾里得長度- 交叉內積 $\mathbf{u}_a^\top \mathbf{u}_b = 0$: 不同的左奇異向量是成對正交的, 線性相關性為零

右奇異向量的正交規範性 (\mathbf{V} 的列) 對於 \mathbf{V} 的任何列向量 \mathbf{v}_a :

$$\mathbf{v}_a^\top \mathbf{v}_a = 1, \quad \mathbf{v}_a^\top \mathbf{v}_b = 0, \quad a \neq b$$

- 自身內積 $\mathbf{v}_a^\top \mathbf{v}_a = 1$: 每個右奇異基底向量都具有單位的歐幾里得長度- 交叉內積 $\mathbf{v}_a^\top \mathbf{v}_b = 0$: 不同的右奇異向量是成對正交的, 線性相關性為零

主成分分析的解讀 這些單位長度和成對正交的條件保證了SVD 基底向量定義了不相關的主成分。零交叉內積消除了奇異方向之間的線性依賴, 這正是用於降維的正交主成分的決定性數學屬性。

數值正交規範性示範例子

我們使用緊湊的4 維矩陣 \mathbf{V} 的兩個列來數值驗證正交規範性規則:

$$\mathbf{v}_1 = \begin{bmatrix} 0.5147 \\ 0.1832 \\ -0.6413 \\ -0.5382 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 0.3102 \\ 0.9221 \\ 0.1658 \\ 0.1739 \end{bmatrix}$$

1. 自身內積（單位長度條件）：

$$\begin{aligned}\mathbf{v}_1^\top \mathbf{v}_1 &= (0.5147)^2 + (0.1832)^2 + (-0.6413)^2 + (-0.5382)^2 \\ &= 0.2649 + 0.0336 + 0.4113 + 0.2897 \\ &= 1.0000\end{aligned}$$

2. 交叉內積（成對正交條件）：

$$\begin{aligned}\mathbf{v}_1^\top \mathbf{v}_2 &= (0.5147)(0.3102) + (0.1832)(0.9221) + (-0.6413)(0.1658) + (-0.5382)(0.1739) \\ &= 0.1597 + 0.1689 - 0.1063 - 0.0936 \\ &\approx 0.0000\end{aligned}$$

同樣的結果適用於 \mathbf{U} 和 \mathbf{V} 中的每一對列。在四位小數四捨五入的容差範圍內，所有自身內積皆等於1，所有交叉內積皆等於0。