

主成分分析

主成分分析(*principal component analysis, PCA*)是一种常用的无监督学习方法,这一方法利用**正交变换**把由**线性相关变量**表示的观测数据转换为少数几个由**线性无关变量**表示的数据,线性无关的变量称为**主成分**。并且主成分的个数通常小于原始变量的个数,所以主成分分析属于降维方法。在数据总体上进行的分析称为总体主成分分析,在有限样本上进行的主成分分析称为样本主成分分析。

基础知识

- 正交矩阵:

如果 n 阶矩阵 A 满足 $A^T A = E$ 即 $A^{-1} = A^T$,那么称 A 为正交矩阵。方阵 A 为正交矩阵的充分必要条件是 A 的列向量都是单位向量且两两正交。

- 方差:

设 x 是一个随机变量,若 $E\{[x - E(x)]^2\}$ 存在,则称其为 x 的方差,记为 $Var(x)$ 或 $D(x)$ 。

- 协方差:

称 $E\{[x - E(x)][y - E(y)]\}$ 为随机变量 x, y 的协方差,记为 $Cov(x, y)$ 。

- 相关系数:

$\rho_{xy} = \frac{Cov(x, y)}{\sqrt{D(x)}\sqrt{D(y)}}$ 称为随机变量 x, y 的相关系数。

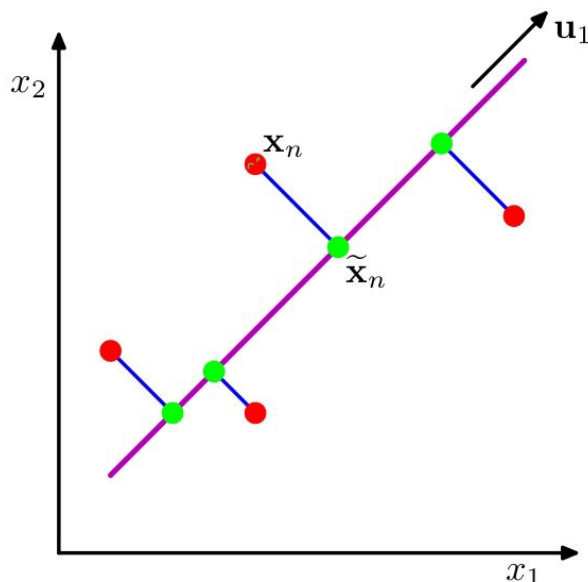
- 协方差矩阵:

- 向量内积的几何意义:

向量 a 和 b 的内积可以表示为 $a \cdot b = |a| \times |b| \cos \theta$,几何意义就是一个向量在另一个向量上的投影与这个向量模长的积,也就是同方向的积。特别地,如果一个向量 a 是某个坐标轴的单位坐标向量,那么向量的内积自然就是 $|b| \cos \theta$,也就是向量 b 在此坐标轴上的坐标值。因此,如果想要将一个向量变换到新的坐标系,那么只要对新坐标系向量进行内积运算即可。

总体主成分分析

基本想法



定义和导出

假设 $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$ 是 m 维随机变量，意即每一维都是单独的随机变量，但是相关性未知。设其均值向量是 $\boldsymbol{\mu} = E(\mathbf{x}) = (\mu_1, \mu_2, \dots, \mu_m)^T$ ，协方差矩阵是 $\boldsymbol{\Sigma} = Cov(\mathbf{x}, \mathbf{x}) = E\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T\}$ 。

$$\begin{aligned}\boldsymbol{\Sigma} &= \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \\ \vdots \\ x_m - \mu_m \end{bmatrix} \begin{bmatrix} x_1 - \mu_1 & x_2 - \mu_2 & \cdots & x_m - \mu_m \end{bmatrix} \\ &= \begin{bmatrix} (x_1 - \mu_1)^2 & (x_1 - \mu_1)(x_2 - \mu_2) & \cdots & (x_1 - \mu_1)(x_m - \mu_m) \\ (x_2 - \mu_2)(x_1 - \mu_1) & (x_2 - \mu_2)^2 & \cdots & (x_2 - \mu_2)(x_m - \mu_m) \\ \vdots & \vdots & \ddots & \vdots \\ (x_m - \mu_m)(x_1 - \mu_1) & (x_m - \mu_m)(x_2 - \mu_2) & \cdots & (x_m - \mu_m)^2 \end{bmatrix}\end{aligned}$$

以上可以看出协方差矩阵 $\boldsymbol{\Sigma}$ 是对称矩阵。

考虑由 m 维随机变量 \mathbf{x} 到 m 维随机变量 $\mathbf{y} = (y_1, y_2, \dots, y_m)$ 的线性变换 $\boldsymbol{\alpha}_i^T$ ：

$\boldsymbol{\alpha}_i^T$ 本质上是个实数向量，不是随机变量。

$$y_i = \boldsymbol{\alpha}_i^T \mathbf{x} = \alpha_{i1}^T x_1 + \alpha_{i2}^T x_2 + \cdots + \alpha_{im}^T x_m$$

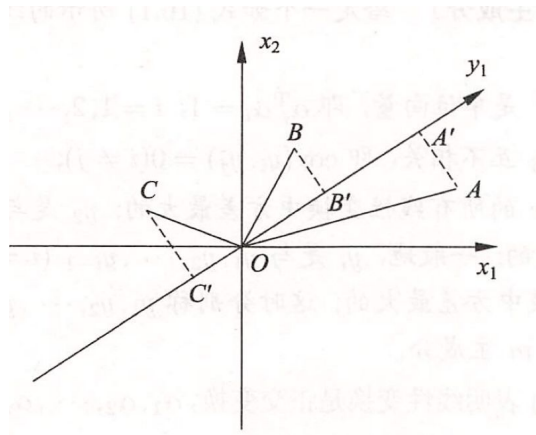
由随机变量的性质可知：

$$E(y_i) = \boldsymbol{\alpha}_i^T E(\mathbf{x}) = \boldsymbol{\alpha}_i^T \boldsymbol{\mu}$$

$$\begin{aligned}Var(y_i) &= E\{[y_i - E(y_i)]^2\} \\ &= E\{(\boldsymbol{\alpha}_i^T \mathbf{x} - \boldsymbol{\alpha}_i^T \boldsymbol{\mu})^2\} \\ &= E\{[\boldsymbol{\alpha}_i^T (\mathbf{x} - \boldsymbol{\mu})][\boldsymbol{\alpha}_i^T (\mathbf{x} - \boldsymbol{\mu})]\} \\ &= E\{[\boldsymbol{\alpha}_i^T (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\alpha}_i]\} \\ &= \boldsymbol{\alpha}_i^T E\{[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T]\} \boldsymbol{\alpha}_i \\ &= \boldsymbol{\alpha}_i^T \boldsymbol{\Sigma} \boldsymbol{\alpha}_i\end{aligned}$$

$$\begin{aligned}Cov(y_i, y_j) &= E\{(y_i - E(y_i))(y_j - E(y_j))\} \\ &= E\{(\boldsymbol{\alpha}_i^T \mathbf{x} - \boldsymbol{\alpha}_i^T \boldsymbol{\mu})(\boldsymbol{\alpha}_j^T \mathbf{x} - \boldsymbol{\alpha}_j^T \boldsymbol{\mu})\} \\ &= E\{(\boldsymbol{\alpha}_i^T \mathbf{x} \mathbf{x}^T \boldsymbol{\alpha}_j - 2\boldsymbol{\alpha}_i^T \boldsymbol{\mu} \mathbf{x} \boldsymbol{\alpha}_j + \boldsymbol{\alpha}_i^T \boldsymbol{\mu} \boldsymbol{\mu}^T \boldsymbol{\alpha}_j)\} \\ &= E\{\boldsymbol{\alpha}_i^T ((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T) \boldsymbol{\alpha}_j\} \\ &= \boldsymbol{\alpha}_i^T \boldsymbol{\Sigma} \boldsymbol{\alpha}_j\end{aligned}$$

值得注意的是：以上三个的最终结果均是一个数。



主成分分析对数据进行正交变换，具体地，对原坐标系进行旋转变换，并将数据在新坐标系表示。比如原数据由变量 (x_1, x_2) 表征，那么变换后在新坐标系里就由变量 (y_1, y_2) 表征。主成分分析选择方差最大的方向作为新坐标系的第一坐标轴 y_1 轴，然后选择与第一坐标轴正交，且方差最大的方向作为新坐标系的第二坐标轴 y_2 轴。

首先，我们来求 x 的第一主成分 $y_1 = \alpha_1^T x$ ，即求系数向量 α_1 ，这个 α_1 实际上就是变换后的坐标轴。因为 α_1 是单位向量，所以由内积的几何意义得知， y_1 就是 x 变换后在 α_1 上的坐标。而选取方差最大的变量，也就是旋转变换后坐标值的平方和最大的轴。

定理：设 x 是 m 维随机变量， Σ 是 x 的协方差矩阵， Σ 的特征值分别为 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$ ，特征值对应的单位特征向量分别为 $\alpha_1, \alpha_2, \dots, \alpha_m$ ，则 x 的第 k 个主成分是

$$y_k = \alpha_k^T x = \alpha_{1k}^T x_1 + \alpha_{2k}^T x_2 + \dots + \alpha_{mk}^T x_m$$

x 的第 k 个主成分的方差为

$$Var(y_k) = \alpha_k^T \Sigma \alpha_k = \lambda_k$$

推论： m 维随机变量 $y = (y_1, y_2, \dots, y_m)^T$ 的分量依次是 x 的第一主成分到第 m 主成分的充要条件是：

- $y = A^T x$ ， A 为正交矩阵。
- y 的协方差矩阵为对角矩阵即 $Cov(y) = diag(\lambda_1, \lambda_2, \dots, \lambda_m)$ ，其中 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$ 。 λ_k 是 x 的协方差矩阵的第 k 个特征值， α_k 是对应的单位特征向量。

总体主成分的性质

- 总体主成分 y 的协方差矩阵是对角矩阵，对角线上的元素即是每个分量的方差。
- 总体主成分 y 的方差之和等于随机变量 x 的方差之和，即 $\sum_{i=1}^m \lambda_i = \sum_{i=1}^m \sigma_{ii}$ 。其中 σ_{ii} 是随机变量 x_i 的方差，即协方差矩阵 Σ 的对角元素。实际上是说， y 的协方差矩阵的迹等于 x 的协方差矩阵的迹。
- 第 k 个主成分与变量 x_i 的相关系数 $\rho(y_k, x_i)$ 称为因子负荷量。计算公式为

$$\rho(y_k, x_i) = \frac{Cov(y_k, x_i)}{\sqrt{Var(y_k)Var(x_i)}} = \frac{\sqrt{\lambda_k} \alpha_{ik}}{\sqrt{\sigma_{ii}}}$$

- 第 k 个主成分 y_k 与 m 个变量的因子负荷量满足 $\sum_{i=1}^m \alpha_{ii} \rho^2(y_k, x_i) = \lambda_k$ 。
- m 个主成分与第 i 个变量 x_i 的因子负荷量满足 $\sum_{i=1}^m \rho^2(y_k, x_i) = 1$ 。

样本主成分分析

上述的主成分分析是定义在样本总体上的，在实际问题中，需要在观测数据上进行主成分分析，这就是样本主成分分析。

定义和性质

假设对 m 维随机变量 $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$ 进行 n 次独立观测， $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ 表示观测样本，其中 $\mathbf{x}_j = (x_{1j}, x_{2j}, \dots, x_{mj})^T$ 表示第 j 个观测样本， x_{ij} 表示第 j 个观测样本的第 i 个变量，则 n 次观测数据可以用样本矩阵 \mathbf{X} 表示，记作：

$$\mathbf{X} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_n] = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}$$

给定了样本我们就可以估计样本均值 $\bar{\mathbf{x}}$ 和样本协方差 \mathbf{S} ：

$$\begin{aligned} \bar{\mathbf{x}} &= \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \\ \mathbf{S} &= [s_{ij}]_{m \times m} \\ s_{ij} &= \frac{1}{n-1} \sum_{k=1}^n (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j) \end{aligned}$$

定义 m 维向量 $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$ 到 m 维向量 $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$ 的线性变换 $\mathbf{y} = \mathbf{A}^T \mathbf{x}$ 。

其中

$$\mathbf{A} = [\mathbf{a}_1 \quad \mathbf{a}_2 \quad \dots \quad \mathbf{a}_m] = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mm} \end{bmatrix}$$

考虑 \mathbf{y} 中的任意的一个线性变换

$$y_i = \mathbf{a}_i^T \mathbf{x} = a_{1i}^T x_1 + a_{2i}^T x_2 + \dots + a_{mi}^T x_m$$

其中 y_i 是 m 维向量 \mathbf{y} 的第 i 个随机变量，均值为 $\bar{y}_i = \frac{1}{n} \sum_{j=1}^n a_i^T x_j = a_i^T \bar{\mathbf{x}}$ 。 y_i 的样本方差

$$\begin{aligned} \text{Var}(y_i) &= \frac{1}{n-1} \sum_{j=1}^n (a_i^T x_j - a_i^T \bar{\mathbf{x}})^2 \\ &= a_i^T \left[\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{\mathbf{x}})(x_j - \bar{\mathbf{x}})^T \right] a_i = a_i^T \mathbf{S} a_i \end{aligned}$$

对任意两个线性变换 $y_i = \mathbf{a}_i^T \mathbf{x}$, $y_k = \mathbf{a}_k^T \mathbf{x}$ ，样本协方差为 $\text{Cov}(y_i, y_k) = a_i^T \mathbf{S} a_k$ 。

对于一组给定的数据点 $\{v_1, v_2, \dots, v_n\}$, 中心化后的表示为

$\{x_1, x_2, \dots, x_n\} = \{v_1 - \mu, v_2 - \mu, \dots, v_n - \mu\}$, 其中 $\mu = \frac{1}{n} \sum_{i=1}^n v_i$ 。由向量内积的几何意义得出, 我们要找一个投影方向 a 使得 x_1, x_2, \dots, x_n 在 a 上的投影坐标方差尽可能大。由于我们提前中心化的原因, 投影之后均值为0, ($\mu' = \frac{1}{n} \sum_{i=1}^n x_i^T a = (\frac{1}{n} \sum_{i=1}^n x_i^T) a = 0$), 这正是提前中心化的意义。

因此投影后的方差可以表示为

$$\begin{aligned} D(x) &= \frac{1}{n} \sum_{i=1}^n (x_i^T a)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (a^T x_i x_i^T a) \\ &= a^T \sum_{i=1}^n \left(\frac{1}{n} x_i x_i^T \right) a \end{aligned}$$