

PCA

主成分分析(Principal Components Analysis, 简称PCA), 用于数据降维

PCA与最大可分性

对于 $X = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix}$, 希望将 X 从 n 维降到 n' 维, 同时希望信息损失最少

基变换

基变换(线性变换): $Y = PX$, 其中 P 是基向量, X 是原始样本, Y 是新表达。

$$\begin{bmatrix} p_1 \\ p_2 \\ \dots \\ p_r \end{bmatrix}_{r \times n} \begin{bmatrix} x_1 & x_2 & \dots & x_m \end{bmatrix}_{n \times m} = \begin{bmatrix} p_1 x_1 & p_1 x_2 & \dots & p_1 x_m \\ \dots & \dots & \dots & \dots \\ p_r x_1 & p_r x_2 & \dots & p_r x_m \end{bmatrix}_{r \times m}$$

其中 p_i 表示行向量, 表示第 i 个基; x_j 是列向量, 表示第 j 个原始数据记录。当基的维度 $r <$ 数据维度 n 时, 可以达到降维的目的

方差

我们希望投影后的数据尽量分散开, 用方差来表达 分散程度, 方差越大, 数据越分散。 $Var(a) = \frac{1}{m} \sum_{i=1}^m (a_i - \mu)^2$, 为了方便处理, 一般将每个字段内所有值都减去字段的平均值

协方差与协方差矩阵

- $Cov(a, b) = \frac{1}{m} \sum_{i=1}^m a_i b_i$

2. 对于 n 维随机变量,

$$x_i = \begin{bmatrix} x_1 \\ x_2 \\ \cdots \\ x_n \end{bmatrix} C = \begin{bmatrix} Var(x_1) & Cov(x_1, x_2) & \cdots & Cov(x_1, x_n) \\ Cov(x_2, x_1) & Var(x_2) & \cdots & Cov(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ Cov(x_n, x_1) & \cdots & \cdots & Var(x_n) \end{bmatrix}$$

即可得到协方

差矩阵。

假设有 m 个 n 维数组记录, 将其按列排成 $n \times m$ 的矩阵 \mathbf{X} , 则 $C = \frac{1}{m} \mathbf{X} \mathbf{X}^T$ 就是协方差矩阵。

3. 协方差矩阵对角化

假设原始数据矩阵为 $X_{n \times m}$, 想要找到基 $P_{r \times n}$ 使得 $Y_{r \times m} = P_{r \times n} X_{n \times m}$ 实现降维的目的。假设 X 的协方差矩阵为 C , Y 的协方差矩阵为 D , 且有 $Y = PX$ 。目标是让协方差矩阵 D 的各个方向方差最大。

则经过化简有: $D = \frac{1}{m} Y Y^T = P C P^T$ 。存在单位向量 $E = [e_1 \ e_2 \ \cdots \ e_n]$,

$$E^T C E = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix}, \text{ 因此 } P = E^T$$

PCA算法过程

设有 m 条 n 维数据

1. 将原始数据按列组成 n 行 m 列矩阵 X
2. 将 X 的每一行(代表一个特征)进行零均值化(减去这一行的均值)
3. 求出协方差矩阵 $C = \frac{1}{m} X X^T$
4. 求出协方差矩阵 C 的特征值以及对应的特征向量
5. 将特征向量按对应特征值大小从上到下按行排列成矩阵, 取前 k 行组成矩阵 P
6. $Y = PX$ 即为降维到 k 维后的矩阵