

## 主成分分析 (PCA)

主成分分析是通过一系列的计算，分析得到样本 $d$  个属性的重要度排序，然后就可以选择性的进行降维了，降到多少维( $d'$ )，就取 $\text{top } d'$  的属性就好了。

---

输入：样本集  $D = \{x_1, x_2, \dots, x_m\}$ ;  
低维空间维数  $d'$ 。

过程：

- 1: 对所有样本进行中心化:  $x_i \leftarrow x_i - \frac{1}{m} \sum_{i=1}^m x_i$ ;
- 2: 计算样本的协方差矩阵  $XX^T$ ;
- 3: 对协方差矩阵  $XX^T$  做特征值分解;
- 4: 取最大的  $d'$  个特征值所对应的特征向量  $w_1, w_2, \dots, w_{d'}$ 。

输出：投影矩阵  $W = (w_1, w_2, \dots, w_{d'})$ 。

---

图 10.5 PCA 算法 <https://blog.csdn.net/pcgamer>

现在假设有一组数据如下：

	x	y
	2.5	2.4
	0.5	0.7
	2.2	2.9
	1.9	2.2
Data =	3.1	3.0
	2.3	2.7
	2	1.6
	1	1.1
	1.5	1.6
	1.1	0.9

行代表了样例，列代表特征，这里有10个样例，每个样例两个特征。可以这样认为，有10篇文档， $x$ 是10篇文档中“learn”出现的TF-IDF， $y$ 是10篇文档中“study”出现的TF-IDF。第一步，分别求 $x$ 和 $y$ 的平均值，然后对于所有的样例，都减去对应的均值。这里 $x$ 的均值是1.81， $y$ 的均值是1.91，那么一个样例减去均值后即为 (0.69,0.49) ，得到

	x	y
	0.69	0.49
	-1.31	-1.21
	0.39	0.99
	0.09	0.29
DataAdjust =	1.29	1.09
	0.49	0.79
	0.19	-0.31
	-0.81	-0.81
	-0.31	-0.31
	-0.71	-1.01

第二步，求特征协方差矩阵，如果数据是3维，那么协方差矩阵是

$$C = \begin{pmatrix} cov(x, x) & cov(x, y) & cov(x, z) \\ cov(y, x) & cov(y, y) & cov(y, z) \\ cov(z, x) & cov(z, y) & cov(z, z) \end{pmatrix}$$

这里只有x和y，求解得

$$cov = \begin{pmatrix} .616555556 & .615444444 \\ .615444444 & .716555556 \end{pmatrix}$$

对角线上分别是x和y的方差，非对角线上是协方差。协方差是衡量两个变量同时变化的变化程度。协方差大于0表示x和y若一个增，另一个也增；小于0表示一个增，一个减。如果x和y是统计独立的，那么二者之间的协方差就是0；但是协方差是0，并不能说明x和y是独立的。协方差绝对值越大，两者对彼此的影响越大，反之越小。协方差是没有单位的量，因此，如果同样的两个变量所采用的量纲发生变化，它们的协方差也会产生树枝上的变化。

第三步，求协方差的特征值和特征向量，得到

$$eigenvalues = \begin{pmatrix} .0490833989 \\ 1.28402771 \end{pmatrix}$$

$$eigenvectors = \begin{pmatrix} -.735178656 & -.677873399 \\ .677873399 & -.735178656 \end{pmatrix}$$

上面是两个特征值，下面是对应的特征向量，特征值0.0490833989对应特征向量为，这里的特征向量都归一化为单位向量。

第四步，将特征值按照从大到小的顺序排序，选择其中最大的k个，然后将其对应的k个特征向量分别作为列向量组成特征向量矩阵。

这里特征值只有两个，我们选择其中最大的那个，这里是1.28402771，对应的特征向量是

$(-0.677873399, -0.735178656)^T$ 。

第五步，将样本点投影到选取的特征向量上。假设样例数为 $m$ ，特征数为 $n$ ，减去均值后的样本矩阵为 $DataAdjust(m \times n)$ ，协方差矩阵是 $n \times n$ ，选取的 $k$ 个特征向量组成的矩阵为 $EigenVectors(n \times k)$ 。那么投影后的数据 $FinalData$ 为

$FinalData(10 \times 1) = DataAdjust(10 \times 2 \text{ 矩阵}) \times \text{特征向量}(-0.677873399, -0.735178656)^T$ 得到的结果是

Transformed Data (Single rigenvector)	
x	
-0.827970186	
1.77758033	
-0.992197494	
-0.274210416	
-1.67580142	
-0.912949103	
0.991094375	
1.14457216	
0.438046137	
1.22382056	

这样，就将原始样例的 $n$ 维特征变成了 $k$ 维，这 $k$ 维就是原始特征在 $k$ 维上的投影。上面的数据可以认为是 $learn$ 和 $study$ 特征融合为一个新的特征叫做 $LS$ 特征，该特征基本上代表了这两个特征。上述过程如下图2描述：

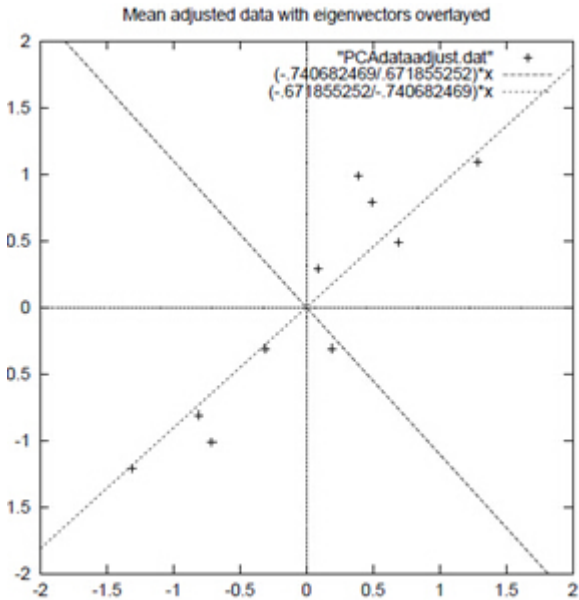


图 2

正号表示预处理后的样本点，斜着的两条线就分别是正交的特征向量（由于协方差矩阵是对称的，因此其特征向量正交），最后一步的矩阵乘法就是将原始样本点分别往特征向量对应的轴上做投影。

整个PCA过程貌似及其简单，就是求协方差的特征值和特征向量，然后做数据转换。但是没有觉得很神奇，为什么求协方差的特征向量就是最理想的k维向量？其背后隐藏的意义是什么？整个PCA的意义是什么？

[后续讲解](#)