

# 主成分分析 (principle component ananalysis)

Guangyao Zhao

2022-05-10

## Contents

样本均值 . . . . .	1
样本方差矩阵 . . . . .	2
最大投影差 . . . . .	2
代码 . . . . .	3

样本表示:

$$X = (x_1; x_2; \dots; x_n) = \begin{pmatrix} x_{11} & x_{11} & \cdots & x_{1m} \\ x_{21} & x_{21} & \cdots & x_{2m} \\ \vdots & & & \\ x_{n1} & x_{n1} & \cdots & x_{nm} \end{pmatrix}$$

其中  $n$  个样本,  $m$  个维度。

## 样本均值

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} (x_1, x_2, \dots, x_n) \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}_n = \frac{1}{n} X^T \mathbf{1}_n$$

其中:

$$\mathbf{1}_n = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}_n$$

## 样本方差矩阵

$$\begin{aligned} S &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T \\ &= \frac{1}{n} (x_1 - \bar{x}, x_2 - \bar{x}, \dots, x_n - \bar{x}) \begin{pmatrix} (x_1 - \bar{x})^T \\ (x_2 - \bar{x})^T \\ \vdots \\ (x_n - \bar{x})^T \end{pmatrix} \\ &= \frac{1}{n} \mathbf{X}^T (\mathbf{I}_n - \bar{\mathbf{X}} \mathbf{1}_n \mathbf{1}_n^T) (\mathbf{X}^T (\mathbf{I}_n - \bar{\mathbf{X}} \mathbf{1}_n \mathbf{1}_n^T))^T \end{aligned}$$

将  $\bar{\mathbf{X}} = \frac{1}{n} \mathbf{X}^T \mathbf{1}_n$  代入上式子:

$$\begin{aligned} S &= \frac{1}{n} \left[ \mathbf{X}^T \left( \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right) \right] \left[ \mathbf{X}^T \left( \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right) \right]^T \\ &= \frac{1}{n} \mathbf{X}^T \mathbf{H} \mathbf{H}^T \mathbf{X} \\ &= \frac{1}{n} \mathbf{X}^T \mathbf{H} \mathbf{X} \end{aligned}$$

其中:  $\mathbf{H} = \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$ 。H 为中心矩阵。

由以上可证明协方差矩阵  $S$  是对称矩阵, 可对角化。

## 最大投影差

- 一个中心: 将一组可能线性相关的变量, 变成一组线性无关的变量。即, 对原始特征的重构。
- 两个基本点: 最大投影方差; 最小重构距离。

向量  $a$  在向量  $u$  方向的投影计算公式为:

$$d = \frac{au}{u}$$

其中  $u$  为单位向量, 即  $u^T u = 1$ 。

PCA 的思想是将原有空间的样本尽可能分散地映射到一个子空间，即各维度应该最大。且尽可能的各维度线性无关，即使基向量正交。根据线性代数，我们可以知道同一元素的协方差就表示该元素的方差，不同元素之间的协方差就表示它们的相关性。

此时就可将问题转化为（首先需要中心化）：

$$\begin{aligned} J &= \sum_{i=1}^n ((x_i - \bar{x})^T u_1)^2, \\ &= \sum_{i=1}^n u_1^T (x_i - \bar{x})(x_i - \bar{x})^T u_1 \\ &= u_1^T \left( \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T \right) u_1 \\ &= u_1^T S u_1 \end{aligned}$$

求  $J$  最大值，将其转化为优化问题：

$$u_1^* = \underset{u_1}{\operatorname{argmax}} u_1^T S u_1$$

上式子满足： $s.t. u^T u = 1$ 。引入拉格朗日乘法：

$$\begin{aligned} L(u_1, \lambda) &= u_1^T S u_1 + \lambda(1 - u_1^T u_1) \\ \frac{\partial L(u_1, \lambda)}{\partial u_1} &= 2S u_1 - 2\lambda u_1 = 0 \end{aligned}$$

即： $S u_1 = \lambda u_1$ 。其中  $\lambda$  为特征值， $u_1$  为特征向量。

以上推导过程说明：信息量保存能力最大的基向量一定是样本矩阵  $X$  的协方差矩阵的特征向量，并且这个特征向量保存的信息量就是它对应的特征值的绝对值。这个推导过程就解释了为什么 PCA 算法要利用样本协方差的特征向量矩阵来降维。

## 代码

```
1  ##Python 实现 PCA
2  import numpy as np
3  import matplotlib.pyplot as plt
4
5  def pca(X, k):
6
```

```

7   m = X.shape[1]
8   norm_X = X - np.mean(X, axis=0) # 中心化
9   scatter_matrix = np.dot(np.transpose(norm_X), norm_X) # 协方差矩阵
10
11  #Calculate the eigenvectors and eigenvalues
12  eig_val, eig_vec = np.linalg.eig(scatter_matrix)
13  eig_pairs = [(np.abs(eig_val[i]), eig_vec[:, i]) for i in range(m)]
14  print('eig_vec:', eig_vec)
15  print('eig_pairs: ', eig_pairs)
16
17  eig_pairs.sort(reverse=True) # 按照特征值从大到小排序
18  feature = np.array([ele[1] for ele in eig_pairs[:k]]) # 选择最大的 k 个特征向量
19
20  #get new data
21  data = np.dot(norm_X, feature.T) # 降维
22
23  # 绘图
24  x = X[:, 0].flatten()
25  y = X[:, 1].flatten()
26  plt.scatter(x, y)
27  plt.plot([eig_vec[:, 0][0], 0], [eig_vec[:, 0][1], 0],
28           color='red') # 特征向量 1
29  plt.plot([eig_vec[:, 1][0], 0], [eig_vec[:, 1][1], 0], color='y') # 特征向量 2
30  plt.show()
31  return data
32
33  X = np.array([[ -1,  1], [-2, -1], [-3, -2], [ 1,  1], [ 2,  1], [ 3,  2]])
34  print(pca(X, 1))

```

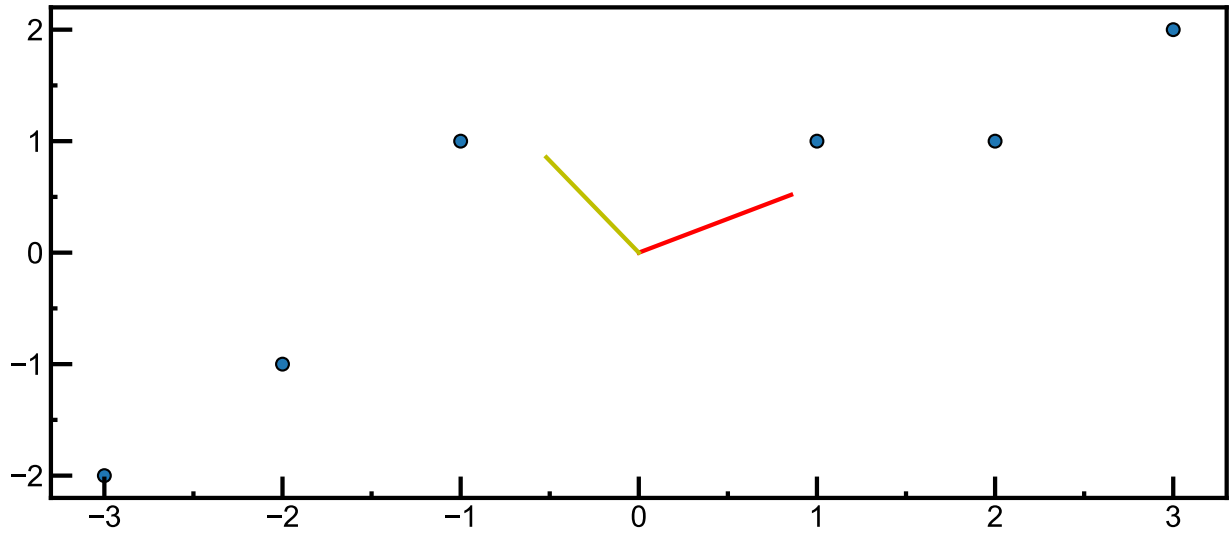
```
eig_vec: [[ 0.8549662 -0.51868371]
```

```
 [ 0.51868371  0.8549662 ]]
```

```
eig_pairs: [(37.70674550364641, array([0.8549662 , 0.51868371])), (1.6265878296869225, array([-0.51868371,  0.8549662 ]))]
```

```
[[ -0.50917706]
```

```
 [-2.40151069]
```



[-3.7751606 ]  
[ 1.20075534]  
[ 2.05572155]  
[ 3.42937146]