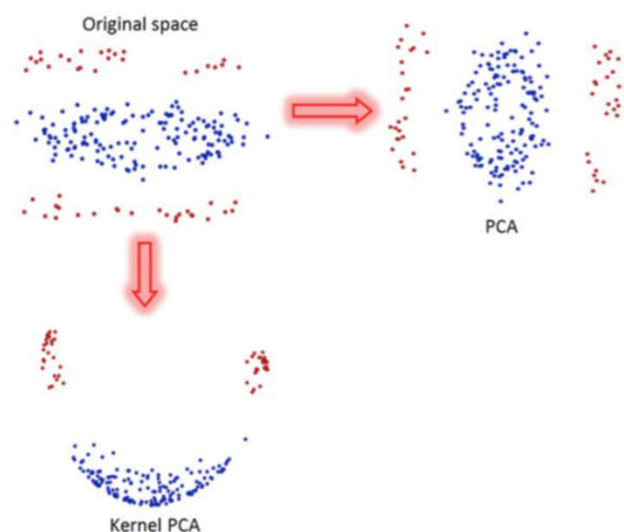


Python编程与人工智能实践

算法篇：数据降维-Kernel-PCA

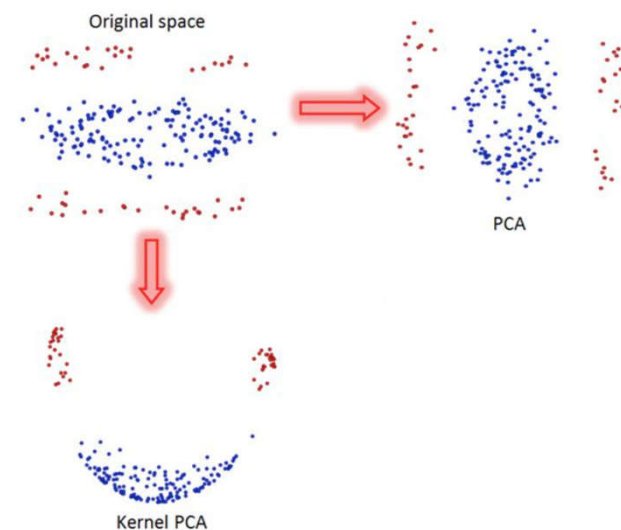


于泓
鲁东大学
信息与电气工程学院
2021.9.27

Kernel-PCA (核PCA)

- PCA降维是一种线性变换，在高维空间的样本点如果线性不可分的话，到了低维空间依旧线性不可分。为了能够让样本在低维空间线性可分引入了KPCA(Kernel PCA)

KPCA基本思想：对源空间的点进行，**非线性高维映射**，在映射后的后的空间内再进行**PCA**降维



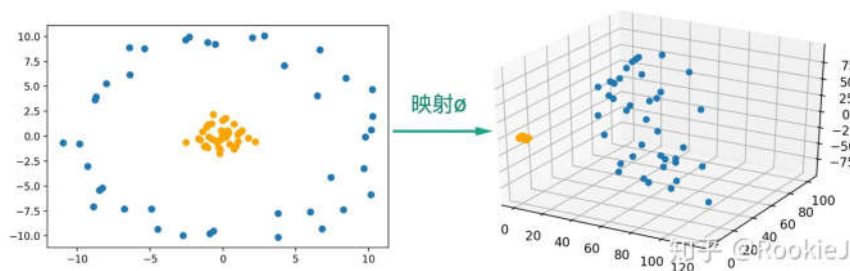
非线性高维映射：（方便表达：后继所有的矢量都是行向量）

例： $\mathbf{x} = [x_1, x_2]$ （低维）

$$\mathbf{X} = \phi(\mathbf{x}) = [x_1^2, x_2^2, \sqrt{2}x_1x_2] \quad (\text{高维})$$

通过映射函数 $\phi(\mathbf{x})$

将一个2维的样本点 \mathbf{x}
映射为3维的样本点 \mathbf{X}



在实际的应用中 $\phi(\mathbf{x})$ 往往未知，已知的通常为核函数

$$k(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{X}_1 \mathbf{X}_2^T = \phi(\mathbf{x}_1) \phi(\mathbf{x}_2)^T$$

即两个高维空间点之间的相关性

常见的核函数：

多项式核

$$k(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{a} \mathbf{x}_1 \mathbf{x}_2^T + c)^d$$

高斯核/RBF(径向基)

$$k(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|^2\right)$$

Sigmoid核

$$k(\mathbf{x}_1, \mathbf{x}_2) = \tanh\left(\mathbf{a} \mathbf{x}_1 \mathbf{x}_2^T + r\right)$$

这些核函数都相当于 $\phi(\mathbf{x})$ 映射的特征维度为无穷多维

一般PCA的步骤:

(1) 对输入的数据 \mathbf{X} (维度 $[N, D]$, N 个样本, 每个 D 维度)

进行中心化处理 $\mathbf{x} = \mathbf{x} - \bar{\mathbf{x}}$

(2) 计算相关矩阵 $\mathbf{C} = \frac{1}{N} \mathbf{X}^T \mathbf{X}$ $\mathbf{C} \in [D, D]$

(特征点各个维度之间的相关性)

(3) 对 \mathbf{C} 进行特征值分解, 获取 d 个较大的特征值, 所对应的特征向量 \mathbf{v} $\mathbf{v} \in [d, D]$

(4) 降维: $\mathbf{x}' = \mathbf{x} \mathbf{v}^T$ $\mathbf{x}' \in [N, d]$

在KPCA中

$$\mathbf{C} = \frac{1}{N} \mathbf{X}^T \mathbf{X} \quad \mathbf{C} \in [D, D]$$

由于映射函数 $\phi(\mathbf{x})$ 通常未知

所以无法求 \mathbf{C} 无法进行特征值分解

已知:

$$\mathbf{K} = \mathbf{X} \mathbf{X}^T \quad \mathbf{K} \in [N, N]$$

各个特征点之间的相关性

设法利用 \mathbf{K} 求 \mathbf{C} 的特征向量

对 $\mathbf{K}=\mathbf{X}\mathbf{X}^T$ $\mathbf{K} \in [N, N]$ 进行特征值分解

$$(\mathbf{X}\mathbf{X}^T)\mathbf{u}^T = \lambda \mathbf{u}^T \quad \mathbf{u} \in [1, N] \quad \mathbf{X} \in [N, D]$$

$$\mathbf{X}^T (\mathbf{X}\mathbf{X}^T)\mathbf{u}^T = \lambda \mathbf{X}^T \mathbf{u}^T$$

$$(\mathbf{X}^T \mathbf{X})(\mathbf{u}\mathbf{X})^T = \lambda (\mathbf{u}\mathbf{X})^T$$

对 $\mathbf{C}=\mathbf{X}^T \mathbf{X}$ $\mathbf{C} \in [D, D]$

进行特征值分解

$$(\mathbf{X}^T \mathbf{X})\mathbf{v}^T = \lambda \mathbf{v}^T \quad \mathbf{v} \in [1, D]$$

得到 \mathbf{C} 的特征向量就是 $\mathbf{v}=\mathbf{u}\mathbf{X}$ $\mathbf{v} \in [1, D]$

对特征矢量 $\mathbf{v} = \mathbf{u}\mathbf{X}$ $\mathbf{v} \in [1, D]$ $\mathbf{X} \in [N, D]$ $\mathbf{u} \in [1, N]$ 进行归一化

$$\begin{aligned} \mathbf{v} &= \frac{1}{\|\mathbf{u}\mathbf{X}\|} \mathbf{u}\mathbf{X} = \frac{1}{\sqrt{\mathbf{u}\mathbf{X}(\mathbf{u}\mathbf{X})^T}} \mathbf{u}\mathbf{X} = \frac{1}{\sqrt{\mathbf{u}(\mathbf{X}\mathbf{X}^T)\mathbf{u}^T}} \mathbf{u}\mathbf{X} & (\mathbf{X}\mathbf{X}^T)\mathbf{u}^T &= \lambda \mathbf{u}^T \\ &= \frac{1}{\sqrt{\mathbf{u}\lambda \mathbf{u}^T}} \mathbf{u}\mathbf{X} = \frac{1}{\sqrt{\lambda}} \mathbf{u}\mathbf{X} = \mathbf{u}'\mathbf{X} \end{aligned}$$

未知

收集特征值较大的 d 个 \mathbf{v} 组成 \mathbf{V} $\mathbf{V} \in [d, D]$

$$\mathbf{V} = \mathbf{U}'\mathbf{X} \quad \mathbf{U}' \in [d, N]$$

降维:

$$\begin{aligned} \mathbf{x}' &= \mathbf{X}\mathbf{v}^T \quad \mathbf{x}' \in [N, d] \\ &= \mathbf{X}\mathbf{X}^T\mathbf{u}'^T = \mathbf{K}\mathbf{u}'^T \end{aligned}$$

$$\mathbf{K} \in [N, N]$$

中心化

$$\tilde{X}_i = X_i - \bar{X}$$

写成矩阵的形式

$$\tilde{\mathbf{X}} = \mathbf{X} - \frac{1}{N} \mathbf{I}_{N \times N} \mathbf{X} = \mathbf{X} - \mathbf{I}_N \mathbf{X}$$

$N \times N$ 全1矩阵

$X_i = \phi(x_i)$ 未知, 转而求 $\tilde{\mathbf{K}}$

$$\tilde{\mathbf{K}} = \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T = (\mathbf{X} - \mathbf{I}_N \mathbf{X})(\mathbf{X} - \mathbf{I}_N \mathbf{X})^T$$

$$= (\mathbf{X} - \mathbf{I}_N \mathbf{X})(\mathbf{X}^T - \mathbf{X}^T \mathbf{I}_N^T) = \mathbf{X} \mathbf{X}^T - \mathbf{I}_N \mathbf{X} \mathbf{X}^T - \mathbf{X} \mathbf{X}^T \mathbf{I}_N^T + \mathbf{I}_N \mathbf{X} \mathbf{X}^T \mathbf{I}_N^T$$

$$= \mathbf{K} - \mathbf{I}_N \mathbf{K} - \mathbf{K} \mathbf{I}_N^T + \mathbf{I}_N \mathbf{K} \mathbf{I}_N^T$$

$[N,D]:[3,2]$

例如:

$$\tilde{\mathbf{X}} = \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \\ X_{31} & X_{32} \end{bmatrix} - \frac{1}{3} \begin{bmatrix} X_{11} + X_{21} + X_{31} & X_{12} + X_{22} + X_{32} \\ X_{11} + X_{21} + X_{31} & X_{12} + X_{22} + X_{32} \\ X_{11} + X_{21} + X_{31} & X_{12} + X_{22} + X_{32} \end{bmatrix}$$

$$= \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \\ X_{31} & X_{32} \end{bmatrix} - \frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \\ X_{31} & X_{32} \end{bmatrix}$$

Kernel-PCA 的一般步骤

(1) 给定样本 \mathbf{x} (维度 $[N, D]$, N 个样本, 每个 D 维)

(2) 利用核函数 计算样本两两之间的关系矩阵

$$\mathbf{K} \quad \mathbf{K} \in [N, N] \quad \text{其中} \quad K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

(3) 对 \mathbf{K} 进行中心化

$$\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{I}_N \mathbf{K} - \mathbf{K} \mathbf{I}_N^T + \mathbf{I}_N \mathbf{K} \mathbf{I}_N^T$$

(4) 对 $\tilde{\mathbf{K}}$ 进行特征值分解

获取较大的 d 个特征值, λ_i 以及对应的特征向量 \mathbf{u}_i

并对 \mathbf{u}_i 进行正则化

$$\mathbf{u}'_i = \frac{1}{\sqrt{\lambda_i}} \mathbf{u}_i$$

(5) 收集 d 个 \mathbf{u}'_i 生成降维矩阵 \mathbf{u}'

$$\mathbf{u}' \in [d, N]$$

进行降维

$$\mathbf{x}' = \mathbf{K} \mathbf{u}'^T \quad \mathbf{x}' \in [N, d]$$

代码实现:

```
import numpy as np
import matplotlib.pyplot as plt

def sigmoid(x1, x2, a = 0.25, r=1):
    x = np.dot(x1, x2)
    return np.tanh(a*x+r)

def linear(x1, x2, a=1, c=0, d=1):
    x = np.dot(x1, x2)
    x = np.power((a*x+c), d)
    return x

def rbf(x1, x2, gamma = 10):
    x = np.dot((x1-x2), (x1-x2))
    x = np.exp(-gamma*x)
    return x
```

多项式核

$$k(x_1, x_2) = (ax_1x_2^T + c)^d$$

高斯核/RBF(径向基)

$$k(x_1, x_2) = \exp\left(-\gamma \|x_1 - x_2\|^2\right)$$

Sigmoid核

$$k(x_1, x_2) = \tanh(ax_1x_2^T + r)$$

```
def kpca(data, n_dims=2, kernel = rbf):
```

```
    N,D = np.shape(data)
```

```
    K = np.zeros([N,N])
```

```
    # 利用核函数计算K
```

```
    for i in range(N):
```

```
        for j in range(N):
```

```
            K[i,j]=kernel(data[i],data[j])
```

```
    # 对K进行中心化
```

```
    one_n = np.ones((N, N)) / N
```

```
    K = K - one_n.dot(K) - K.dot(one_n) + one_n.dot(K).dot(one_n) →  $\tilde{K} = K - \mathbf{I}_N \mathbf{K} - \mathbf{K} \mathbf{I}_N^T + \mathbf{I}_N \mathbf{K} \mathbf{I}_N^T$ 
```

```
    #计算特征值和特征向量
```

```
    eig_values, eig_vector = np.linalg.eig(K)
```

```
    idx = np.argsort(-eig_values)[:n_dims] # 从大到小排序
```

```
    # 选取较大的特征值
```

```
    eigval = eig_values[idx]
```

```
    eigvector = eig_vector[:, idx] # [N,d]
```

```
    # 进行正则
```

```
    eigval = eigval**(1/2)
```

```
    u = eigvector/eigval.reshape(-1,n_dims) # u [N,d]
```

$$\mathbf{u}'_i = \frac{1}{\sqrt{\lambda_i}} \mathbf{u}_i$$

```
    # 进行降维
```

```
    data_n = np.dot(K, u)
```

```
    return data_n
```

```
def draw_pic(datas, labs):
    plt.cla()
    unique_labs = np.unique(labs)
    colors = [plt.cm.Spectral(each)
               for each in np.linspace(0, 1, len(unique_labs))]
    p=[]
    legends = []
    for i in range(len(unique_labs)):
        index = np.where(labs==unique_labs[i])
        pi = plt.scatter(datas[index, 0], datas[index, 1], c=[colors[i]] )
        p.append(pi)
        legends.append(unique_labs[i])

    plt.legend(p, legends)
    plt.show()

if __name__ == "__main__":
    # 加载数据
    data = np.loadtxt("iris.data", dtype="str", delimiter=',')
    feas = data[:, :-1]
    feas = np.float32(feas)
    labs = data[:, -1]
    kpca(feas)

    # 进行降维
    data_2d = kpca(feas, n_dims=2, kernel=rbf)

    # 绘图
    draw_pic(data_2d, labs)
```

