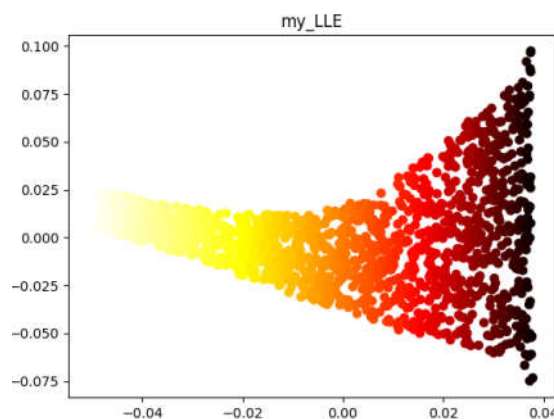
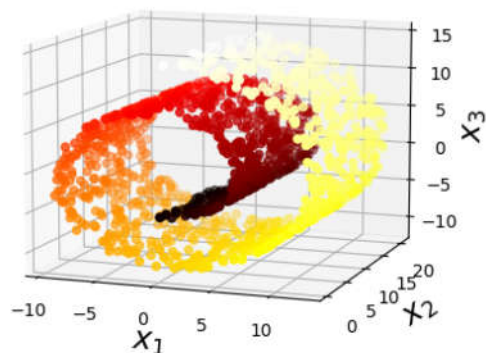


Python编程与人工智能实践

算法篇：数据降维-LLE (Locally Linear Embedding)

局部线性嵌入



于泓
鲁东大学
信息与电气工程学院
2021.10.6

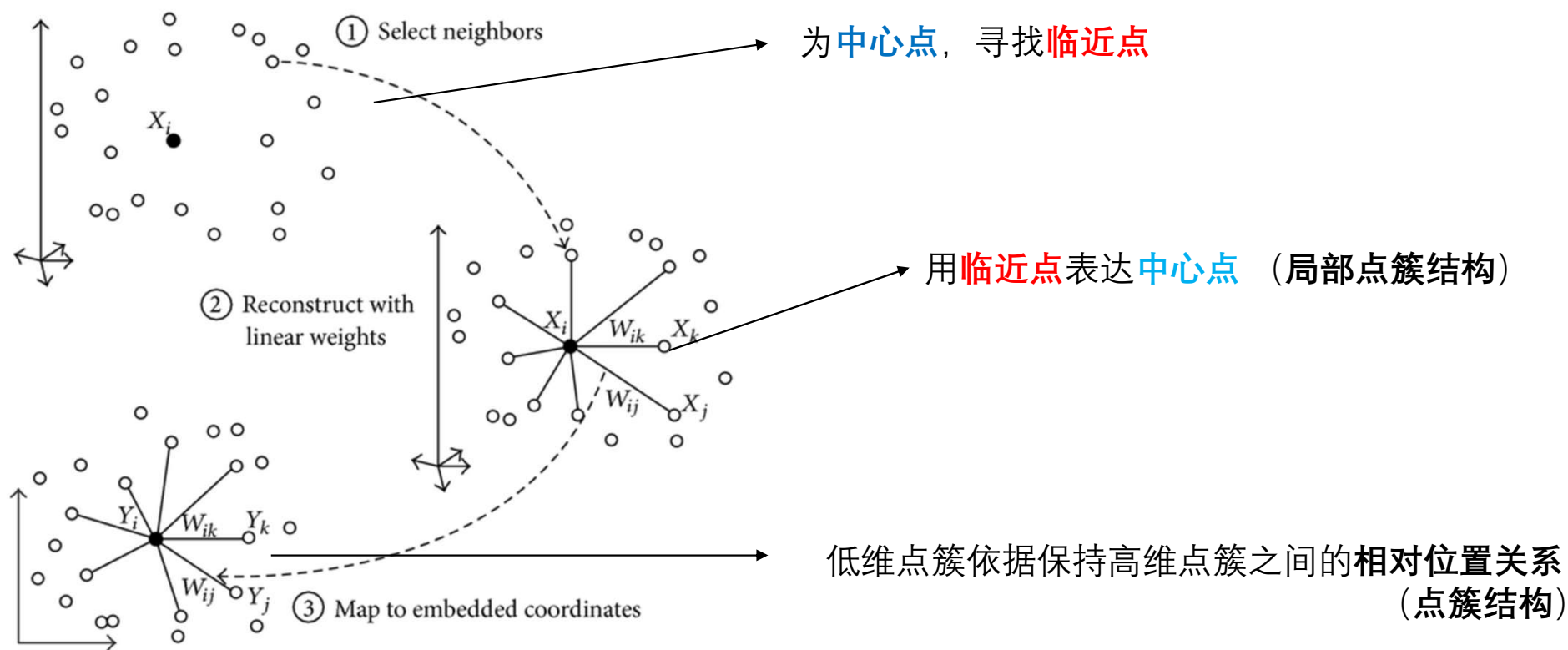
LLE (Locally Linear Embedding, 局部线性嵌入)

LLE (局部线性嵌入) 即“Locally Linear Embedding”的降维算法，属于“流形” (manifold) 降维。在低维空间可以表示高维空间的形状

本法主要原理：在高维空间中，局部点之间的相对位置关系，在低维空间依旧可以保持

对比：MSD算法：保留点与点之间的关系 (考量单独点)

LLE 算法：保留局部分点簇之间的相对关系 (考量局部多个点)



步骤 (1) 高维样本点 X_i 由它的 k 个临近点
进行近似 其中上标 (k) 表示 k 个临近点

$$\arg \min (w_i) = \frac{1}{2} \|X_i - w_i X_i^{(k)}\|^2$$

其中 $X_i \in [1, D]$ $w_i \in [1, k]$ $X_i^{(k)} \in [k, D]$

且要求 $\sum_{j=1}^k w_{ij} = 1$ 即 $w_i \mathbf{I}_{k \times 1} = 1$

$$\begin{aligned} \|X_i - w_i X_i^{(k)}\|^2 &= \|w_i \mathbf{I}_{k \times 1} X_i - w_i X_i^{(k)}\|^2 \\ &= \|w_i (\mathbf{I}_{k \times 1} X_i - X_i^{(k)})\|^2 \\ &= w_i (\mathbf{I}_{k \times 1} X_i - X_i^{(k)}) (\mathbf{I}_{k \times 1} X_i - X_i^{(k)})^T w_i^T = w_i \mathbf{S}_i w_i^T \end{aligned}$$

构建拉格朗日方程

$$L(w_i, \lambda) = \frac{1}{2} w_i \mathbf{S}_i w_i^T - \lambda (w_i \mathbf{I}_{k \times 1} - 1)$$

求导可得 $\frac{\partial L(w_i)}{\partial w_i} = 0$

$$w_i \mathbf{S}_i = \lambda \mathbf{I}_{k \times 1}^T$$

$$w_i = \lambda \mathbf{I}_{k \times 1}^T \mathbf{S}^{-1}$$

$$w_i \mathbf{I}_{k \times 1} = \lambda \mathbf{I}_{k \times 1}^T \mathbf{S}^{-1} \mathbf{I}_{k \times 1} = 1 \rightarrow \lambda = \frac{1}{\mathbf{I}_{k \times 1}^T \mathbf{S}^{-1} \mathbf{I}_{k \times 1}}$$

带入

$$w_i = \frac{\mathbf{I}_{k \times 1}^T \mathbf{S}^{-1}}{\mathbf{I}_{k \times 1}^T \mathbf{S}^{-1} \mathbf{I}_{k \times 1}}$$

列元素的和
 全部元素的和

步骤二 进行低维空间映射

映射后的低维空间点 \mathbf{Y}_i 能够由 k 个临近点 $\mathbf{Y}_i^{(k)}$ \rightarrow k 个临近点的编号
(index 不变)利用步骤一 学的的权重 w_i 进行重构

$$\arg(\mathbf{Y}) = \sum_{i=1}^N \|\mathbf{Y}_i - w_i \mathbf{Y}_i^{(k)}\|^2 \quad \mathbf{Y} \in [N, d]$$

min

约束条件：低维特征，每个维度均值为0，方差为1

可以通过偏置
调节

$$\frac{1}{N} \mathbf{Y}^T \mathbf{Y} = \mathbf{I}_{d \times d}$$

单位阵

$$\sum_{i=1}^N \|\mathbf{Y}_i - w_i \mathbf{Y}_i^{(k)}\|^2 = \sum_{i=1}^N \|\mathbf{Y}_i - \mathbf{W}_i \mathbf{Y}\|^2$$

 $\mathbf{W}_i \in [1, N]$ 其中 k 个近邻元素的位置由 w_i
中的元素填充，其他位置为0

上式可以进一步写为：

$$\arg(\mathbf{Y}) = \underset{\min}{\text{trace}} \left((\mathbf{Y} - \mathbf{WY})(\mathbf{Y} - \mathbf{WY})^T \right)$$

$$= \text{trace} \left(\mathbf{Y}(\mathbf{I} - \mathbf{W})(\mathbf{I} - \mathbf{W})^T \mathbf{Y}^T \right)$$

结合限制条件：构建拉格朗日函数

$$L(\mathbf{Y}, \lambda) = \text{trace} \left(\mathbf{Y}(\mathbf{I} - \mathbf{W})(\mathbf{I} - \mathbf{W})^T \mathbf{Y}^T \right) - \lambda \left(\mathbf{Y}^T \mathbf{Y} - \frac{1}{N} \mathbf{I}_{d \times d} \right)$$

$$L(\mathbf{Y}, \lambda) = \text{trace} \left(\mathbf{Y} (\mathbf{I}_{N \times N} - \mathbf{W}) (\mathbf{I}_{N \times N} - \mathbf{W})^T \mathbf{Y}^T \right) - \lambda \left(\mathbf{Y}^T \mathbf{Y} - \frac{1}{N} \mathbf{I}_{d \times d} \right)$$

求导可得

$$\mathbf{Y} (\mathbf{I}_{N \times N} - \mathbf{W})^T (\mathbf{I}_{N \times N} - \mathbf{W}) = \lambda \mathbf{Y}$$

可以发现降维后的特征 \mathbf{Y}

就是矩阵 $(\mathbf{I}_{N \times N} - \mathbf{W})^T (\mathbf{I}_{N \times N} - \mathbf{W})$ 的特征向量

由于要求最小值，所以应当选取 d 个最小的特征值所对应的特征向量

最小的特征值往往接近 **0**，所以实际应用中，选取 **除最小特征值外** 最小的 d 个特征值所对应的特征向量

作为降维后的数据

LLE算法的一般步骤

(1) 为高维数据 \mathbf{X} $\mathbf{X} \in [N, D]$ 中的每个样本

\mathbf{X}_i 选取 k 个临近点 $\mathbf{X}_i^{(k)}$ $\mathbf{X}_i^{(k)} \in [k, D]$

(2) 为每个样本点计算一组权重 $\mathbf{w}_i \in [1, k]$

可以借助这组权重, 用 $\mathbf{X}_i^{(k)}$ 重构 \mathbf{X}_i

$$\mathbf{w}_i = \frac{\mathbf{I}_{k \times 1}^T \mathbf{S}_i^{-1}}{\mathbf{I}_{k \times 1}^T \mathbf{S}_i^{-1} \mathbf{I}_{k \times 1}}$$

$$\mathbf{S}_i = (\mathbf{I}_{k \times 1} \mathbf{X}_i - \mathbf{X}_i^{(k)}) (\mathbf{I}_{k \times 1} \mathbf{X}_i - \mathbf{X}_i^{(k)})^T$$

(3) 将 \mathbf{w}_i $\mathbf{w}_i \in [1, k]$ 扩充为 \mathbf{W}_i $\mathbf{W}_i \in [1, N]$

k 个临近点的位置, $W_{ik} = w_{ik}$ 其他位置为 0

收集 N 个 \mathbf{W}_i 组成 \mathbf{W}

$$\mathbf{C} = (\mathbf{I}_{N \times N} - \mathbf{W})^T (\mathbf{I}_{N \times N} - \mathbf{W})$$

对 \mathbf{C} 进行特征值分解

选取**除最小特征值外**, 最小的 d 个特征值
所对应的特征向量 \mathbf{Y} $\mathbf{Y} \in [N, d]$

就是降维后的特征点。

代码实现

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import make_s_curve
from sklearn.manifold import Isomap
from tqdm import tqdm
from sklearn.manifold import LocallyLinearEmbedding

# x 维度 [N,D]
def cal_pairwise_dist(x):

    N,D = np.shape(x)

    dist = np.zeros([N,N])

    for i in range(N):
        for j in range(N):
            dist[i,j] = np.sqrt(np.dot((x[i]-x[j]),(x[i]-x[j]).T))

    #返回任意两个点之间距离
    return dist
```

```
# 获取每个样本点的 n_neighbors个临近点的位置
def get_n_neighbors(data, n_neighbors = 10):
    dist = cal_pairwise_dist(data)
    dist[dist < 0] = 0
    N = dist.shape[0]
    Index = np.argsort(dist,axis=1)[:,:n_neighbors+1]
    return Index

def scatter_3d(X, y):
    fig = plt.figure()
    ax = plt.axes(projection='3d')
    ax.scatter(X[:, 0], X[:, 1], X[:, 2], c=y, cmap=plt.cm.hot)
    ax.view_init(10, -70)
    ax.set_xlabel("$x_1$", fontsize=18)
    ax.set_ylabel("$x_2$", fontsize=18)
    ax.set_zlabel("$x_3$", fontsize=18)
    plt.show(block=False)
```



```

# data : N,D
def lle(data, n_dims = 2, n_neighbors = 10):
    N,D = np.shape(data)
    if n_neighbors > D:
        tol = 1e-3
    else:
        tol = 0
    # 获取 n_neighbors个临界点的位置
    Index_NN = get_n_neighbors(data,n_neighbors)

    # 计算重构权重
    w = np.zeros([N,n_neighbors])
    for i in range(N):

        X_k = data[Index_NN[i]] # [k,D]
        X_i = [data[i]] # [1,D]
        I = np.ones([n_neighbors,1])

        Si = np.dot((np.dot(I,X_i)-X_k), (np.dot(I,X_i)-X_k).T)

        # 为防止对角线元素过小
        Si = Si+np.eye(n_neighbors)*tol*np.trace(Si)

        Si_inv = np.linalg.pinv(Si)
        w[i] = np.dot(I.T,Si_inv)/(np.dot(np.dot(I.T,Si_inv),I))

```

$$W_i = \frac{I_{k \times 1}^T S_i^{-1}}{I_{k \times 1}^T S_i^{-1} I_{k \times 1}}$$

$$S_i = (I_{k \times 1} X_i - X_i^{(k)})(I_{k \times 1} X_i - X_i^{(k)})^T$$

```

# 计算 W
W = np.zeros([N,N])
for i in range(N):
    W[i,Index_NN[i]] = w[i]

```

```

I_N = np.eye(N)
C = np.dot((I_N-W).T, (I_N-W))

```

$C = (I_{N \times N} - W)^T (I_{N \times N} - W)$

```

# 进行特征值的分解
eig_val, eig_vector = np.linalg.eig(C)

```

```

index_ = np.argsort(eig_val)[1:n_dims+1]

```

```

y = eig_vector[:,index_]
return y

```

实验结果

