

## §10 降维与度量学习

### 1. KNN

机制：给定测试样本，基于某种距离度量找出训练集中与其最靠近的  $k$  个训练样本，然后基于  $k$  个邻居进行预测

分类：投票法 -  $k$  个样本中出现最多的标记 (多数表决, 多数服从少数)  
 回归：平均法 -  $k$  个样本的取值输出标记的平均值作为预测结果

给定测试样本  $x$ , 若最近  $k$  个样本为  $z$ , 则最近邻分类器输出  $z$  的标签, 即

$$P(\text{class}) = \frac{1}{k} \sum_{i=1}^k P(c|x) P(c|z)$$

三要素:  $k$  值的选择 ( $k$  小易过拟合,  $k$  大易欠拟合), 距离度量 (欧氏)  
 决策规则 (分类, 回归)

算法:  $\begin{cases} \text{暴力 brute} \\ kD\text{-tree, ball-tree, BBF tree, MUP tree} \end{cases}$

$kD\text{-tree}$ : 从  $m$  个样本的  $n$  维特征中, 分别计算  $n$  个特征取值的方差, 因方差最大的第  $k$  维特征  $h_k$  作为根节点, 取值  $h_{kl}$  作为划分点,  $<$  该值的为左子树,  $>$  该值的为右子树

### 2. 低维嵌入和 PCA

降维: 通过某种数学变换将原始高维属性空间转换为低维子空间  
 即高维空间中的低维嵌入 多维缩放 (MDS)

PCA: 通过某种线性投影, 将高维的数据映射到低维的空间中表示, 并且期望在所映射的维数上数据的信息量最大, 以此使用较少维数。

正交属性平面中存在着平面。

最近量线性, 样本到该平面是最近

最大可分性, 样本点在这个平面上的投影点可分离

原理:  $X$  是已中心化 (z-score) 过的数据矩阵, 将点  $x_i$  在超平面上投影到  $W^T x_i$ , 投影后点在各轴向上的方差最大化, 方差和  $\frac{1}{n} \sum W^T x_i x_i^T W$ .

目标函数:  $\max_W \text{tr}(W^T X X^T W)$

$$L = W^T X X^T W + \lambda (I - W^T W) \quad \frac{\partial L}{\partial W} = 2 X X^T W - 2 \lambda W$$

$$\text{令 } \frac{\partial L}{\partial W} = 0 \rightarrow X X^T W = \lambda W$$

$$W^T X X^T W = W^T \lambda W = \lambda$$

步骤: 数据中心化  $\rightarrow$  求中心化矩阵  $X$  的协方差矩阵  $\rightarrow$  特征值/向量

$\rightarrow$  从大到小排列, 选取前面的  $k$  个特征向量  $W$

3. 核化线性降维和流形学习

基于核技巧对线性降维进行“核化”

$$k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

$$z_j = W_j^T \phi(x) = \sum_{i=1}^m \alpha_i^j \phi(x_i)^T \phi(x)$$

$$= \sum_{i=1}^m \alpha_i^j k(x, x_i)$$

KPCA 计算量大

流形学习  $\left\{ \begin{array}{l} \text{高维空间映射} \\ \text{局部线性嵌入} \end{array} \right.$  最短距离 MDS 算法  
LLE