

KNN原理小结

KNN分类与回归除了最后一个一般是多数表决策法，一个一般是选择平均法。scikit-learn只使用了蛮力法 (brute force)
KD树实现 (KD-Tree) 和球树 (ball tree)

① KNN算法的三要素.

三个最终的要素: k 值的选取, 距离度量的方式和分类决策规则。[分类决策规则一般用多数表决策法, 重点关注 k 值的选择和距离的度量方式。

k 值的选取, 没有固定的经验, 一般根据样本的维, 选择一个较小的值, 可以通过交叉验证来选择一个合适的 k 值。 k 值太小, 相似的训练实例对结果起作用, 泛化误差大, k 值小 \Rightarrow 模型复杂, 容易发生过拟合。

k 值太大, 整体模型变得简单。泛化误差减小, 训练误差增大。

距离的度量: ①常用的是欧氏距离: 对于两个 n 维向量 $x(x_1, \dots, x_n)$ 和 $y(y_1, \dots, y_n)$ $D(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$

②曼哈顿距离 $|x_1 - y_1| + |x_2 - y_2| + \dots + |x_n - y_n| = \sum_{i=1}^n |x_i - y_i| = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$

③通用闵可夫斯基距离 $\sqrt[p]{|x_1 - y_1|^p + |x_2 - y_2|^p + \dots + |x_n - y_n|^p} = \sqrt[p]{\sum_{i=1}^n (x_i - y_i)^p}$

② KNN的算法蛮力实现.

计算预测样本与训练集中样本的距离, 然后计算最小的 k 个距离即可, 接着多数表法。在特征数多, 样本量大的时候, 算法时间效率很成问题

② KNN算法之KD树实现原理.

KD树没有一开始就对测试样本分类, 而是先对训练集建模, 建之的模型就是KD树。建好模型对 test data 做预测
这里的模型就是KD树: KD树中的 k 代表样本的特征维度。不是KNN中最远的 k 个最近样本的意思。后面用 n 表示

KD tree 三步: 建树、搜索最近邻、预测。

KD树建树: KD树是二叉树, 每个节点有 [特征坐标, 切割轴, 指向左枝的指针, 指向右枝的指针]

切割轴 $r \quad 1 \leq r \leq n$ (n 为特征维度的数量)

对每一个维度计算方差, 方差大的那一维为 r , 取第 r 维的中位数, 假如是 $x_r^{(u)}$ 计 $x_r^{(u)} < x_r^{(u)}$ 分到左边, 对 $x_r^{(u)} > x_r^{(u)}$ 分到右边, (若有 m 个样本, m 为偶数, 则取中位左边或右边。

搜索最近邻: 对测试集 P , 把测试集 P 放在叶子节点上, 持续向上爬, 与到其父节点和切割轴的距离比
挺复杂的, 具体算法见 聚宽 - kd树详解篇。