当样本数少于特征数时,容易过拟合,最好先用PCA, LDA等提取主成分

如果类别分布非常不均匀,就要考虑用class_weight

如果样本中有string,可以用labelEncoder、 OneHotEncoder或者DictVectorizer

如果想把连续型数据类别化,可以用分桶的思想,或者二 值化

基于树的计算都不需要对数据做标准化处理

criterion: gini/entropy 一般选gini

splitter:样本不大选best,反之random 全局最优和局部 最优的关系

max_features: None考虑全部特征,其它的为考虑N的 对数、平方根,看样本大小

max_depth: 一般不超过100

min_samples_split: 一般10万数据控制在10左右

min_samples_leaf: 一般10万数据在5个左右

min_weight_fraction_leaf: 如果想给某个叶子节点特定的权重

max_leaf_nodes: 最大叶子节点数

class_weight: 类别权重,balanced 样本少的类别会赋予 跟高的权重

min_impurity_split:不纯度小于这个阈值,则该节点不再生成子节点。即为叶子节点。

presor:数据是否预排序,一般false

决策树可以认为是一种if-then的集合,即互斥且穷尽。

可以看做是一种条件概率分布,基于特征的类划分。

3个步骤:特征选择、决策树生成、减枝

sklearn.tree.DecisionTreeClassifier参数

总概:决策树DT是一种分类和回归的方

法。

注意

利用信息増益来选择特征 利用信息増益很可能会选有更多属性的特征 sklearn.tree.DecisionTreeClassifier(entropy)

「信息増益比是除法运算 (H(D)-H(D/A))/ splitinformation (D,A)

利用信息増益比来选择特征 惩罚了有更多属性的特征,弥补ID3 sklearn.tree.DecisionTreeClassifier(entropy)

GINI指数越大,样本的不确定性越大 対回归用MSE,分类用GINI sklearn.tree.DecisionTreeClassifier(gini)

信息增益是减法计算,H(D)-H(D/A)

信息增益最大~信息增益比最大~基尼指数最小