

Python机器学习

线性模型

逻辑回归

逻辑回归和线性回归都是线性的，但是分别解决分类问题和回归问题
S型曲线Sigmoid (0, 1) 之间的概率值, $y = \text{logistic}(b + w_1 * x_1 + \dots + w_n * x_n) = 1 / (1 + \exp(-(b + w_1 * x_1 + \dots + w_n * x_n)))$
模型复杂度越高越容易过拟合，加入正则项防止过拟合，平衡损失函数与模型复杂度: $\text{minimize}(\text{Loss}(\text{Data} | \text{Model}) + \text{complexity}(\text{Model}))$ 惩罚特别大的，sklearn中，LR的参数C越大正则化越弱，反之越强

线性回归

$Y = w_1 * x_1 + w_2 * x_2 + \dots + b$ ，用最小二乘法求参数，根据真实值和预测值计算残差平方和，sklearn.linear_model

SVM

$f(x, w, b) = \text{sign}(w * x + b) = \text{sign}(\sum w_i * x_i + b)$ 最初是做二分类，结果是负数或正数，核心就是找最大间隔的分类器，C值越大正则化越弱间隔越窄，反之正则化越强间隔越宽
建树过程：1、从根结点开始，计算所有特征值的信息增益（信息增益比），选择计算结果最大的特征作为根结点
2、根据算出的特征建立子节点，执行第一步，直到所有特征的信息增益很小或没有特征可以选择为止

决策树

防止过拟合：预剪枝和后剪枝，改变树的深度

KNN(K-近邻)

KNeighborsClassifier

fit()

基本概念

监督学习：带标记

作用：预估下个季度的销售量

时序分析

作用：预测业务增长或房价等，需要有标签的

回归预测

线性回归

KNN

分类

逻辑回归

SVM

非监督学习：未被标记

应用：根据症状归纳特定疾病等

聚类

基本流程

1. 数据的表示（特征工程）：定义一系列函数

2. 模型检验：定义函数的优劣

3. 模型优化：选择最优函数

机器学习库scikit-learn

1. 准备数据集

2. 选择模型

3. 训练模型调整参数

4. 测试模型