

## Naive Bayes:

1. 设  $x = \{a_1, a_2, \dots, a_m\}$  为待分类项, 而每个  $a$  为  $x$  的一个特征属性
2. 类别集合  $C = \{y_1, \dots, y_n\}$
3. 算  $p(y_1|x), p(y_2|x), \dots, p(y_n|x)$
4. 若  $p(y_k|x) = \max \{p(y_1|x), \dots, p(y_n|x)\}$ , 则  $x \in y_k$

① 训练集:

② 求  $p(a_1|y_1), p(a_2|y_1), \dots, p(a_m|y_1); p(a_1|y_2), \dots, p(a_m|y_2), \dots; p(a_1|y_n), \dots, p(a_m|y_n)$

③ 若  $a_1, a_2, \dots, a_m$  是独立的, 则  $p(y_i|x) = \frac{p(x|y_i)p(y_i)}{p(x)}$

$\therefore p(x)$  对 all 类别为常数.  $\therefore$  只需  $p(x|y_i)p(y_i)$  最大化

$$p(x|y_i) = p(a_1|y_i) \dots p(a_m|y_i) \cdot p(y_i) = p(y_i) \cdot \prod_{j=1}^m p(a_j|y_i)$$

**NB:** 基础思想: 对于待分类项, 求解 ~~由该~~  $x$  出现的条件下, 各个类别  $y_1, \dots, y_n$  出现的  $x_i$  概率, 取最大概率的类别.

pros: ① 简单 ② 分类准确率较高, 速度快, 所需估计参数少. 对缺失数据不敏感

cons: ① 假设各个属性互不相关, 不可能. ② 需知先验 prob

## Decision Tree

基础思想: 通过训练数据构建决策树, 每个节点, 表示在某个属性上的测试, 每个分枝代表该测试的一个输出, 每个叶节点存放一个类别

ID3: info gain. ; C4.5: 信息增益比; CART: 基尼系数

pros: ① 无需参数假设 ② 适合非线性数据. ③ 简单

cons: ① 易过拟合 ② 忽略属性间相关性 ③ Imbalanced data, info gain 偏向于多样本特征



## bagging 随机森林

特点: 对 outliers 更鲁棒。all 依赖个体分类能力, 和它们之间依赖性, 对每次划分属性很敏感。 (不希望有相关性)

优点: ① 不易 overfit, ② 每次划分只取部分 attr,  $\therefore$  在大型 DB 很有效

③ 即便 missing 很多, 也能维持 high accuracy

④ 对 imbalanced data, 可以平衡误差

缺点: ① 在某些噪声较大的分类、回归问题上会过拟合。

② 级联划分较多的属性会对 RF 产生更大影响

boosting: 对弱分类器封装结合成强分类器。 (串行)

优点: all 高, 无太多参数  
缺点: 对 outlier 敏感

random forest 比一般决策树稳定原因:

bagging, 多树投票提高泛化能力且引入随机, 避免单棵树过拟合, 提高整体泛化能力

## xgboost 与 gbd

gbd 以 CART 作为基分类器 ① xgboost 还支持线性分类器。

② 有  $L_1$   $L_2$  的 regularization 控制模型复杂度

③ 对 loss function 进行 2 阶 Taylor 展开

④ 支持列抽样

⑤ 自动学习 missing value 的分布

⑥ 在特征粒度上并行

⑦ 在训练前, 对 data 排序, 保存为 block 结构, 后面迭代反复用这个结构。

训练时一阶导信息