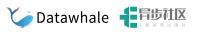


《机器学习公式详解》 (南瓜书)

第7章 贝叶斯分类器

本节主讲: 谢文睿

本节大纲

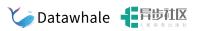


西瓜书对应章节: 7.1、7.2、7.3

- 1. 贝叶斯决策论
- 2. 生成式模型和判别式模型
- 3. 朴素贝叶斯分类器
- 4. 半朴素贝叶斯分类器



贝叶斯决策论是概率框架下实施决策的基本方法,对分类任务来说,在所有相关概率都已知的理想情形下,贝叶斯决策论考虑如何基于这些概率和误判损失来选择最优的类别标记。



以一个多分类任务为例:假设当前有一个N分类问题,即 $\mathcal{Y} = \{c_1, c_2, \ldots, c_N\}$

【定义】: λ_{ij} 是将一个真实标记为 c_j 的样本误分类为 c_i 所产生的损失。

【定义】:单个样本x的期望损失(条件风险)为

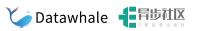
$$R\left(c_{i}|oldsymbol{x}
ight)=\sum_{j=1}^{N}\lambda_{ij}P\left(c_{j}|oldsymbol{x}
ight)$$

其中, $P(c_j|\boldsymbol{x})$ 为后验概率

【定义】: 全部样本构成的总体风险为

$$R(h) = \mathbb{E}_{m{x}}[R(h(m{x})|m{x})]$$

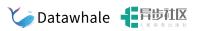
其中,h为分类器(模型)。显然,分类效果越准确的h,其条件风险和总体风险也越小



贝叶斯判定准则:为最小化总体风险R(h),只需在每个样本上选择那个能使条件风险 $R(c|\boldsymbol{x})$ 最小的类别标记,即

$$h^*(oldsymbol{x}) = rg \min_{c \in \mathcal{Y}} \!\! R(c|oldsymbol{x})$$

此时, h^* 称为贝叶斯最优分类器



具体地,若目标是最小化分类错误率,则误判损失 λ_{ij} 可写为

$$\lambda_{ij} = \left\{ egin{array}{ll} 0, & ext{if } i=j \ 1, & ext{otherwise} \end{array}
ight.$$

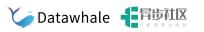
此时单个样本x的期望损失(条件风险)为

$$R\left(c_{i}|oldsymbol{x}
ight)=\sum_{j=1}^{N}\lambda_{ij}P\left(c_{j}|oldsymbol{x}
ight)$$

$$R(c_i|\boldsymbol{x}) = 1 * P(c_1|\boldsymbol{x}) + ... + 1 * P(c_{i-1}|\boldsymbol{x}) + 0 * P(c_i|\boldsymbol{x}) + 1 * P(c_{i+1}|\boldsymbol{x}) + ... + 1 * P(c_N|\boldsymbol{x})$$

又
$$\sum_{j=1}^N P(c_j|oldsymbol{x})=1$$
,则

$$R(c_i|oldsymbol{x}) = 1 - P(c_i|oldsymbol{x})$$

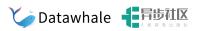


于是,按照贝叶斯判定准则,最小化分类错误率的贝叶斯最优分类器为

$$egin{aligned} h^*(oldsymbol{x}) &= rg \min_{c \in \mathcal{Y}} R(c|oldsymbol{x}) \ h^*(oldsymbol{x}) &= rg \min_{i \in \{1,2,...,N\}} R(c_i|oldsymbol{x}) \ h^*(oldsymbol{x}) &= rg \min_{i \in \{1,2,...,N\}} 1 - P(c_i|oldsymbol{x}) \ h^*(oldsymbol{x}) &= rg \max_{i \in \{1,2,...,N\}} P(c_i|oldsymbol{x}) \ h^*(oldsymbol{x}) &= rg \max_{c \in \mathcal{Y}} P(c|oldsymbol{x}) \end{aligned}$$

即对每个样本 \boldsymbol{x} ,选择后验概率 $P(c_i|\boldsymbol{x})$ 最大的类别 c_i 作为标记

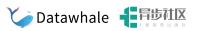
生成式模型和判别式模型



从贝叶斯决策论(概率框架)的角度:机器学习所要做的就是基于有限的训练样本集尽可能准确地估计出后验概率 $P(c|\boldsymbol{x})$

从机器学习自己的角度:给定一个样本 $m{x}$,求一个能准确分类 $m{x}$ 的 $m{f}(m{x})$,其有些算法可以看作是对后验概率建模 $m{P}(c|m{x})$ (例如对数几率回归),而有些算法则是纯粹完成样本分类(例如SVM)

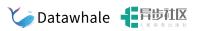
生成式模型和判别式模型



判别式模型: 给定 $oldsymbol{x}$, 直接建模 $P(c|oldsymbol{x})$ 来预测c

生成式模型: 先对联合概率 $P(\boldsymbol{x},c)$ 建模, 然后再由此推导得出 $P(c|\boldsymbol{x})$

生成式模型和判别式模型



对于生成式模型, 其建模思路为

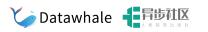
$$P(c|oldsymbol{x}) = rac{P(oldsymbol{x},c)}{P(oldsymbol{x})}$$

再根据贝叶斯定理, 上式可恒等变形为

$$P(c|oldsymbol{x}) = rac{P(c)P(oldsymbol{x}|c)}{P(oldsymbol{x})}$$

其中,P(c)是类"先验"概率, $P(\boldsymbol{x}|c)$ 是样本 \boldsymbol{x} 相对于类别标记c的类条件概率, $P(\boldsymbol{x})$ 是用于归一化的"证据"因子。

朴素贝叶斯分类器



属性条件独立性假设:对已知类别,假设所有属性相互独立

$$P(c|oldsymbol{x}) = rac{P(c)P(oldsymbol{x}|c)}{P(oldsymbol{x})} = rac{P(c)}{P(oldsymbol{x})} \prod_{i=1}^d P\left(x_i|c
ight)$$

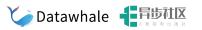
其中,d为属性数目, x_i 为x在第i个属性上的取值。基于贝叶斯判定准则

$$h^*(oldsymbol{x}) = rgmax_{c \in \mathcal{Y}} P(c|oldsymbol{x}) = rgmax_{c \in \mathcal{Y}} rac{P(c)}{P(oldsymbol{x})} \prod_{i=1}^{a} P\left(x_i|c
ight)$$

由于对所有类别来说P(x)都相同,所以P(x)视作常量可以略去

$$h_{nb}(oldsymbol{x}) = rgmax_{c \in \mathcal{Y}} P(c) \prod_{i=1}^d P\left(x_i|c
ight)$$

朴素贝叶斯分类器

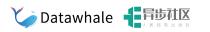


估计P(c):

$$P(c) = rac{|D_c|}{|D|}$$

其中, D_c 表示训练集D中类别标记为c的样本集合, $|D_c|$ 表示集合 D_c 的样本总数

朴素贝叶斯分类器



估计 $P(x_i|c)$:

【第i个属性为离散属性】:

$$P\left(x_i|c
ight) = rac{|D_{c,x_i}|}{|D_c|}$$

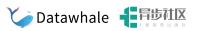
其中, D_{c,x_i} 表示 D_c 中在第i个属性上取值为 x_i 的样本组成的集合

【第i个属性为连续属性】(以正态分布假设为例):

$$p\left(x_{i}|c
ight) = rac{1}{\sqrt{2\pi}\sigma_{c,i}}\exp\left(-rac{\left(x_{i}-\mu_{c,i}
ight)^{2}}{2\sigma_{c,i}^{2}}
ight)$$

其中, $\mu_{c,i}$ 和 $\sigma_{c,i}^2$ 分别是第c类样本在第i个属性上取值的均值和方差。

半朴素贝叶斯分类器



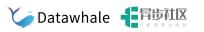
半朴素贝叶斯分类器:适当考虑一部分属性间的相互依赖信息,从而既不需进行完全联合概率计算,又不至于彻底忽略了比较强的属性依赖关系。

【独依赖估计(ODE)】:假设每个属性在类别之外最多依赖于一个其他属性,即

$$P(c|oldsymbol{x}) \propto P(c) \prod_{i=1}^d P\left(x_i|c,pa_i
ight)$$

其中, pa_i 为属性 x_i 所依赖的属性,称为 x_i 的父属性。

半朴素贝叶斯分类器

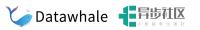


【超父独依赖估计(SPODE)】: 假设所有属性都依赖于同一个"超父"属性

$$egin{aligned} P(c \mid oldsymbol{x}) &= rac{P(oldsymbol{x}, c)}{P(oldsymbol{x})} = rac{P\left(c, x_i
ight) P\left(x_1, \ldots, x_{i-1}, x_{i+1}, \ldots, x_d \mid c, x_i
ight)}{P(oldsymbol{x})} \ &\propto P\left(c, x_i
ight) P\left(x_1, \ldots, x_{i-1}, x_{i+1}, \ldots, x_d \mid c, x_i
ight) \ &= P\left(c, x_i
ight) \prod_{j=1}^d P\left(x_j \mid c, x_i
ight) \end{aligned}$$

其中, x_i 是"超父"属性

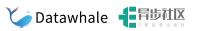
预告



下一节:集成学习

西瓜书对应章节:第8章

结束语



欢迎加入【南瓜书读者交流群】,我们将在群里进行答疑、勘误、本次直播回放、本次直播PPT发放、下次直播通知等最新资源发放和活动通知。加入步骤:

- 1. 关注公众号【Datawhale】,发送【南瓜书】三个字获取机器人"小豚"的微信二维码
- 2. 添加"小豚"为微信好友,然后对"小豚"发送【南瓜书】三个字即可自动邀请进群

