

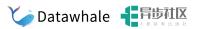


《机器学习公式详解》 (南瓜书)

第8章 集成学习

本节主讲:秦州

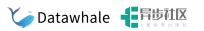
本节大纲



西瓜书对应章节: 8.1、8.2

- 1. 个体与集成
- 2. Adaboost算法

个体与集成

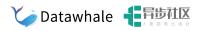


"三个臭皮匠,顶个诸葛亮",集成学习通过集合 **多个个体学习器** 的结果来提升预测结果的准确性和泛化能力。

"君子和而不同"个体学习器需要比随机猜想要强一些,个体学习器的预测结果也要具有一定的多样性。

	样本a	样本b	样本c	样本a	样本b	样本c	样本a	样本b	样本c
学习器1	1	1	0	1	0	0	1	1	0
学习器2	1	0	1	0	1	0	1	1	0
学习器3	0	1	1	0	0	1	1	1	0
集成结果	1	1	1	0	0	0	1	1	0

个体与集成

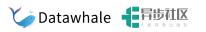


集成个体学习器的收敛性保证:

$$egin{aligned} P(H(oldsymbol{x})
eq \sum_{k=0}^{\lfloor T/2
floor} \left(egin{array}{c} T \ k \end{array}
ight) (1-\epsilon)^k \epsilon^{T-k} \ &\leqslant \exp\left(-rac{1}{2}T(1-2\epsilon)^2
ight) \end{aligned}$$

两个基本结论:

- 收敛速率随着个体学习器数量T呈指数下降
- $\epsilon = 0.5$ 的个体集成器对收敛没有作用



学习T个个体学习器 h_t 和相应的权重 α_t ,使得他们的加权和

$$H(oldsymbol{x}) = \sum_{t=1}^T lpha_t h_t(oldsymbol{x})$$

能够最小化损失函数

$$\ell_{ ext{exp}}(H \mid \mathcal{D}) = \mathbb{E}_{oldsymbol{x} \sim \mathcal{D}} \left[e^{-f(oldsymbol{x})H(oldsymbol{x})}
ight]$$

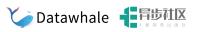


前向分布求解算法: 每一轮只学习一个学习器 h_t 和相应的权重 α_t , 第t轮的优化目标

$$(lpha_t, h_t) = rg\min_{lpha, h} \ell_{ ext{exp}} \left(H_{t-1} + lpha h \mid \mathcal{D}
ight)$$

根据指数损失函数的定义式(8.5),有

$$egin{aligned} \ell_{ ext{exp}}\left(H_{t-1} + lpha h \mid \mathcal{D}
ight) &= \mathbb{E}_{oldsymbol{x} \sim \mathcal{D}}\left[e^{-f(oldsymbol{x})(H_{t-1}(oldsymbol{x}) + lpha h(oldsymbol{x}))}
ight] \ &= \sum_{i=1}^{|D|} \mathcal{D}\left(oldsymbol{x}_i
ight) e^{-f(oldsymbol{x}_i)H_{t-1}(oldsymbol{x}_i)} e^{-f(oldsymbol{x}_i)lpha h(oldsymbol{x}_i)} \ &= \sum_{i=1}^{|D|} \mathcal{D}\left(oldsymbol{x}_i
ight) e^{-f(oldsymbol{x}_i)H_{t-1}(oldsymbol{x}_i)} e^{-f(oldsymbol{x}_i)lpha h(oldsymbol{x}_i)} \end{aligned}$$

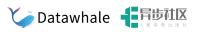


因为 $f(x_i)$ 和 $h(x_i)$ 仅可取值 $\{-1,1\}$,可以推得

$$egin{aligned} \ell_{ ext{exp}}\left(H_{t-1} + lpha h \mid \mathcal{D}
ight) &= \sum_{i=1}^{|D|} \mathcal{D}\left(oldsymbol{x}_i
ight) e^{-f\left(oldsymbol{x}_i
ight)H_{t-1}\left(oldsymbol{x}_i
ight)}\left(e^{-lpha} + \left(e^{lpha} - e^{-lpha}
ight) \mathbb{I}\left(f\left(oldsymbol{x}_i
ight)
eq h\left(oldsymbol{x}_i
ight)
ight) \ &= \sum_{i=1}^{|D|} \mathcal{D}\left(oldsymbol{x}_i
ight) e^{-f\left(oldsymbol{x}_i
ight)H_{t-1}\left(oldsymbol{x}_i
ight)} e^{-lpha} + \sum_{i=1}^{|D|} \mathcal{D}\left(oldsymbol{x}_i
ight) e^{-f\left(oldsymbol{x}_i
ight)H_{t-1}\left(oldsymbol{x}_i
ight)}\left(e^{lpha} - e^{-lpha}
ight) \mathbb{I}\left(f\left(oldsymbol{x}_i
ight)
eq h\left(oldsymbol{x}_i
ight) \end{aligned}$$

做一个简单的符号替换,令 $\mathcal{D}_t'\left(\boldsymbol{x}_i\right)=\mathcal{D}\left(\boldsymbol{x}_i\right)e^{-f\left(\boldsymbol{x}_i\right)H_{t-1}\left(\boldsymbol{x}_i\right)}$,并且注意到 $e^{-\alpha}$ 和 $e^{\alpha}-e^{-\alpha}$ 与求和变量i无关,可以提取出来,有

$$\ell_{ ext{exp}}\left(H_{t-1} + lpha h \mid \mathcal{D}
ight) = e^{-lpha} \sum_{i=1}^{|D|} \mathcal{D}_t'\left(oldsymbol{x}_i
ight) + \left(e^{lpha} - e^{-lpha}
ight) \sum_{i=1}^{|D|} \mathcal{D}_t'\left(oldsymbol{x}_i
ight) \mathbb{I}\left(f\left(oldsymbol{x}_i
ight)
eq h\left(oldsymbol{x}_i
ight)
ight)$$



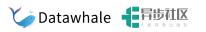
$$\ell_{ ext{exp}}\left(H_{t-1} + lpha h \mid \mathcal{D}
ight) = e^{-lpha} \sum_{i=1}^{|D|} \mathcal{D}_t'\left(oldsymbol{x}_i
ight) + \left(e^{lpha} - e^{-lpha}
ight) \sum_{i=1}^{|D|} \mathcal{D}_t'\left(oldsymbol{x}_i
ight) \mathbb{I}\left(f\left(oldsymbol{x}_i
ight)
eq h\left(oldsymbol{x}_i
ight)
ight)$$

我们的目的是求解 h_t 使得 $\ell_{\rm exp}$ 最小化,因此可以忽略掉与h无关的项,即求解目标是

$$h_t = rg\min_h \left(e^{lpha} - e^{-lpha}
ight) \sum_{i=1}^{|D|} \mathcal{D}_t'\left(oldsymbol{x}_i
ight) \mathbb{I}\left(f\left(oldsymbol{x}_i
ight)
eq h\left(oldsymbol{x}_i
ight)
ight)$$

更近一步,由于 $\alpha>\frac{1}{2}$,易证得 $e^{\alpha}-e^{-\alpha}>0$ 恒成立,因此求解目标为:

$$h_t = rg\min_h \sum_{i=1}^{|D|} \mathcal{D}_t'\left(oldsymbol{x}_i
ight) \mathbb{I}\left(f\left(oldsymbol{x}_i
ight)
eq h\left(oldsymbol{x}_i
ight)
ight)$$



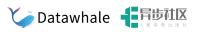
$$h_t = rg\min_h \sum_{i=1}^{|D|} \mathcal{D}_t'\left(oldsymbol{x}_i
ight) \mathbb{I}\left(f\left(oldsymbol{x}_i
ight)
eq h\left(oldsymbol{x}_i
ight)
ight)$$

其中 $\mathcal{D}_t'\left(oldsymbol{x}_i
ight) = \mathcal{D}\left(oldsymbol{x}_i
ight) e^{-f\left(oldsymbol{x}_i
ight)H_{t-1}\left(oldsymbol{x}_i
ight)}$

观察 $\mathcal{D}'_t(\boldsymbol{x}_i)$ 的形式可以发现它仅与t-1轮及以前的学习器有关,因此在求解 h_t 时,对于每个样本i,他其实已经固定了,如果把 $\mathcal{D}'_t(\boldsymbol{x}_i)$ 看做样本i在t轮学习时的权重分布,我们要依据这个权重求解学习器 h_t 以满足上面的最优化式子。

同时,为了确保 $\mathcal{D}_t'(\boldsymbol{x}_i)$ 是一个分布,通常我们对其进行规范化后作为下一个学习器的输入样本权重,即 $\mathcal{D}_t(\boldsymbol{x}_i) = \frac{\mathcal{D}_t'(\boldsymbol{x}_i)}{\sum_{i=1}^{|D|} \mathcal{D}_t'(\boldsymbol{x}_i)}$,其中分母是常数,因此这个变换不会影响上述

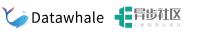
最小化的求解。



有意思的一点是,t轮的样本权重可以通过t-1轮样本权重计算,而无需从头算起,以t+1轮为例,根据迭代公式,有:

$$egin{aligned} \mathcal{D}_{t+1}\left(oldsymbol{x}_i
ight) &= \mathcal{D}\left(oldsymbol{x}_i
ight) e^{-f(oldsymbol{x}_i)H_t(oldsymbol{x}_i)} \ &= \mathcal{D}\left(oldsymbol{x}_i
ight) e^{-f(oldsymbol{x}_i)(H_{t-1}(oldsymbol{x}_i)+lpha_t h_t(oldsymbol{x}_i))} \ &= \mathcal{D}\left(oldsymbol{x}_i
ight) e^{-f(oldsymbol{x}_i)H_{t-1}(oldsymbol{x}_i)} e^{-f(oldsymbol{x}_i)lpha_t h_t(oldsymbol{x}_i)} \ &= \mathcal{D}_t\left(oldsymbol{x}_i
ight) e^{-f(oldsymbol{x}_i)lpha_t h_t(oldsymbol{x}_i)} \end{aligned}$$

这便是《机器学习》式8.19



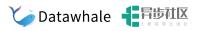
下面求解学习器 h_t 的权重 α_t 。损失函数 $\ell_{\mathrm{exp}}\left(H_{t-1}+\alpha h\mid \mathcal{D}\right)$ 对lpha求导有:

$$\frac{\partial \ell_{\exp}\left(H_{t-1} + \alpha h_t \mid \mathcal{D}\right)}{\partial \alpha} = \frac{\partial \left(e^{-\alpha} \sum_{i=1}^{|D|} \mathcal{D}_t'\left(\boldsymbol{x}_i\right) + \left(e^{\alpha} - e^{-\alpha}\right) \sum_{i=1}^{|D|} \mathcal{D}_t'\left(\boldsymbol{x}_i\right) \mathbb{I}\left(f\left(\boldsymbol{x}_i\right) \neq h\left(\boldsymbol{x}_i\right)\right)\right)}{\partial \alpha}$$

$$= -e^{-\alpha} \sum_{i=1}^{|D|} \mathcal{D}_t'\left(\boldsymbol{x}_i\right) + \left(e^{\alpha} + e^{-\alpha}\right) \sum_{i=1}^{|D|} \mathcal{D}_t'\left(\boldsymbol{x}_i\right) \mathbb{I}\left(f\left(\boldsymbol{x}_i\right) \neq h\left(\boldsymbol{x}_i\right)\right)$$

令导数等于0,移项可得:

$$egin{aligned} rac{e^{-lpha}}{e^{lpha}+e^{-lpha}} &= rac{\sum_{i=1}^{|D|} \mathcal{D}_t'\left(oldsymbol{x}_i
ight) \mathbb{I}\left(f\left(oldsymbol{x}_i
ight)
eq h\left(oldsymbol{x}_i
ight)
ight)}{\sum_{i=1}^{|D|} \mathcal{D}_t'\left(oldsymbol{x}_i
ight)} = \sum_{i=1}^{|D|} rac{\mathcal{D}_t'\left(oldsymbol{x}_i
ight)}{Z_t} \mathbb{I}\left(f\left(oldsymbol{x}_i
ight)
eq h\left(oldsymbol{x}_i
ight)
ight) \\ &= \sum_{i=1}^{|D|} \mathcal{D}_t\left(oldsymbol{x}_i
ight) \mathbb{I}\left(f\left(oldsymbol{x}_i
ight)
eq h\left(oldsymbol{x}_i
ight)
ight) = \mathbb{E}_{oldsymbol{x}\sim\mathcal{D}_t}\left[\mathbb{I}\left(f\left(oldsymbol{x}_i
ight)
eq h\left(oldsymbol{x}_i
ight)
ight)
ight] = \epsilon_t \end{aligned}$$



$$\frac{e^{-\alpha}}{e^{\alpha} + e^{-\alpha}} = \frac{1}{e^{2\alpha} + 1} \Rightarrow e^{2\alpha} + 1 = \frac{1}{\epsilon_t} \Rightarrow e^{2\alpha} = \frac{1 - \epsilon_t}{\epsilon_t} \Rightarrow 2\alpha = \ln\left(\frac{1 - \epsilon_t}{\epsilon_t}\right)$$
$$\Rightarrow \alpha_t = \frac{1}{2}\ln\left(\frac{1 - \epsilon_t}{\epsilon_t}\right)$$

这便是《机器学习》式8.11

当 $\epsilon > \frac{1}{2}$ 时,上式单调递减,因此误差率越大的学习器分配的权重越少。

Datawhale 手步在区

初始化样本权值分布. 基于分布 \mathcal{D}_t 从数据集 D 中训练出分类器 h_t .

估计 ht 的误差.

确定分类器 ht 的权重.

更新样本分布, 其中 Z_t 是规范化因子, 以确保 D_{t+1} 是一个分布.

输入: 训练集
$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\};$$

基学习算法 $\mathfrak{L};$
训练轮数 $T.$
过程:
1: $\mathcal{D}_1(x) = 1/m.$
2: for $t = 1, 2, \dots, T$ do
3: $h_t = \mathfrak{L}(D, \mathcal{D}_t);$
4: $\epsilon_t = P_{x \sim \mathcal{D}_t}(h_t(x) \neq f(x));$
5: if $\epsilon_t > 0.5$ then break
6: $\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right);$

8: end for

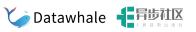
输出:
$$H(x) = \operatorname{sign}\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$$

图 8.3 AdaBoost算法

 $\mathcal{D}_{t+1}(\boldsymbol{x}) = \frac{\mathcal{D}_{t}(\boldsymbol{x})}{Z_t} imes \left\{ egin{array}{ll} \exp(-lpha_t), & ext{if } h_t(\boldsymbol{x}) = f(\boldsymbol{x}) \ \exp(lpha_t), & ext{if } h_t(\boldsymbol{x})
eq f(\boldsymbol{x}) \end{array}
ight.$

 $= \frac{\mathcal{D}_t(\boldsymbol{x}) \exp(-\alpha_t f(\boldsymbol{x}) h_t(\boldsymbol{x}))}{Z_t}$

预告



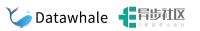
下一节: Bagging、随机森林、多样性增强

增补知识点: GB(Gradient Boosting)/GBDT/XGBoost

西瓜书对应章节: 8.3 8.5

增补知识点参考资料:略

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

