

VC 学习笔记-Hoeffding 不等式

赵得涛 15220172202792

2019 年 10 月 25 日

1 背景

VC 维的概念最早是在1971年, V. Vapnik and A. Chervonenkis在论文“On the uniform convergence of relative frequencies of events to their probabilities”中提出的。是机器学习中重要的概念。在机器学习中, 我们不知道真实模型, 我们只能用近似模型估计真实模型。我们的目标是 $E_{in} \approx 0$ 和 $E_{in} \approx E_{out}$, 通常为了使 E_{in} 最小, 人们会将使用比较复杂的模型(VC很大), 甚至过度拟合。然而这样只是满足的条件一却使条件二偏离很大。VC 维反映了函数集的学习性能, VC 维越大, 学习机器的学习能力越强, 但学习机器也越复杂。

2 Hoeffding 不等式

2.1 定义

Hoeffding不等式是关于一组随机变量均值的概率不等式。如果 X_1, X_2, \dots, X_n 为一组独立同分布的参数为 p 的伯努利分布随机变量, n 为随机变量的个数, 定义这组随机变量的均值为:

$$\bar{X} = \frac{X_1 + X_2 + X_3 + \dots + X_n}{n}$$

对于任意的 $\delta > 0$, Hoeffding不等式可以表示为:

$$P(|\bar{X} - E(\bar{X})| \geq \delta) \leq \exp(-2\delta^2 n)$$

如果知道 $X_i \in [a_i, b_i]$ 时, Hoeffding的应用更加广泛:

$$P(\bar{X} - E(\bar{X}) \geq \delta) \leq \exp - \frac{2n^2 \delta^2}{\sum_{i=1}^n (b_i - a_i)^2}$$

$$P(|\bar{X} - E(\bar{X})| \geq \delta) \leq 2 \exp - \frac{2n^2 \delta^2}{\sum_{i=1}^n (b_i - a_i)^2}$$

从不等式的证明可以看出当 n 趋于无限大时, 我们可以用 \bar{X} 有效推断 $E(\bar{X})$

2.2 Learning 中的应用

假设一个学习算法的实际假设是 $f(x)$ (一般是未知的), $h(x)$ 是采样之后的假设, y_n 为对应样本的目标值: 对于任意固定的 h :

$$E_{out}(h) = \epsilon[h(x) \neq f(x)]$$

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N [h(x) \neq f(x)]$$

根据Hoeffding不等式, 我们可以得出:

$$P(\|E_{in} - E_{out}\| > \epsilon) \leq 2\exp(-2\epsilon^2 N)$$

对任意的 $\epsilon > 0$ 成立, N 表示样本大小, 当 N 越大, 说明采样假设对于真实假设的推断是越好的。Hoeffding不等式只能保证 $E_{in}(h)$ 和 $E_{out}(h)$ 会有一个较大的几率说明他俩是相似的, 但是无法保证在非常多得假设下 ($h(x)$) 下, $E_{in}(h)$ 的值时比较小, 那如果最终 $E_{in}(h)$ 的值比较大, 那么就背离了我们学习的两个目标之一。

假定假设空间 H 中有 M 个假设 $h_1, h_2, h_3, \dots, h_M$, 从输入空间以分布 P 抽取样本数量为 N 的训练集 $D_1, D_2, D_3 \dots D_{5678} \dots$ (可以是无穷个), 如下图:

	\mathcal{D}_1	\mathcal{D}_2	...	\mathcal{D}_{1126}	...	\mathcal{D}_{5678}	Hoeffding
h_1	BAD					BAD	$\mathbb{P}_D[\text{BAD } \mathcal{D} \text{ for } h_1] \leq \dots$
h_2		BAD					$\mathbb{P}_D[\text{BAD } \mathcal{D} \text{ for } h_2] \leq \dots$
h_3	BAD	BAD				BAD	$\mathbb{P}_D[\text{BAD } \mathcal{D} \text{ for } h_3] \leq \dots$
...							
h_M	BAD					BAD	$\mathbb{P}_D[\text{BAD } \mathcal{D} \text{ for } h_M] \leq \dots$
all	BAD	BAD				BAD	?

图 1: 矩阵表示

BAD表示当前的假设下对应的样本 E_{in} 很大。对于任何一个 h , 任何一个训练集, 发生BAD的概率最大为 $2\exp(-2\epsilon^2 N)$, 对于 M 个假设, 我们可以使用bound的方式来求总体的BAD几率:

$$\begin{aligned}
 P(\text{BAD}|D) &= P_D[\text{BAD for } h_1 \text{ or } \text{BAD for } h_2 \text{ or } \dots \text{ or } \text{BAD for } h_M] \\
 &\leq P_D[\text{BAD for } h_1] + P_D[\text{BAD for } h_2] + \dots + P_D[\text{BAD for } h_M] \\
 &\leq 2\exp(-2\epsilon^2 N) + 2\exp(-2\epsilon^2 N) + \dots + 2\exp(-2\epsilon^2 N) \\
 &= 2M\exp(-2\epsilon^2 N)
 \end{aligned}$$

2.3 关于M无限大的问题

拿最简单线性模型来说, 假设空间 $H = \{\beta_0 + \beta_1 X\}$, 我们可以有无数个 β_0 和 β_1 估计, 这样是否意味着概率会变成无穷? 首先开始几个PLA(线性可分)的例子:

当平面上只有一个线性可分的点时, 有两条线(一虚一实), 如下图: 当

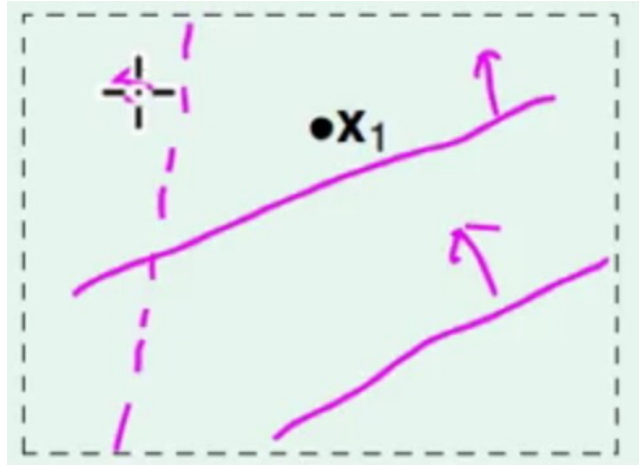


图 2: 1个线性可分得点

平面上只有一个线性可分的点时, 有四条线, 如下图: 在二维上, 平面上

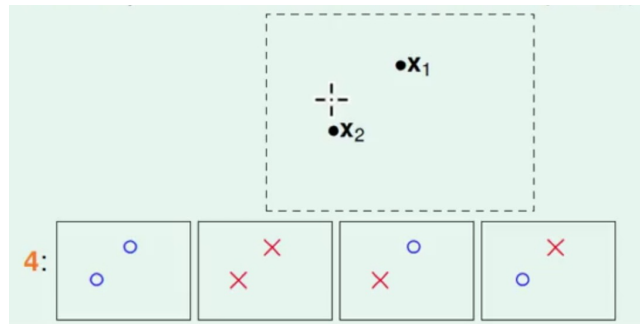


图 3: 2个线性可分得点

划分这些点的种类数总是小于 2^n , 现在假设使用 $M_H(N)$ 来表示 N 个点的时候可能出现的最大分类线种数 dichotomies:

$$M_H(N) = \max \|H(X_1, X_2, \dots, X_N)\|$$

那么这个 $M_H(N)$ 的上界就是 2^N , 这个函数也叫做成长函数 (Growth Function). 如果 $M_H(N)$ 能直接取代 M , 同时 $M_H(N)$ 为多项式的话, N 越大, 那个

这个upbound会越来越小，或者接近于0，如果 $M_H(N)$ 为指数的话，则不是很妙：

$$P(\|E_{in} - E_{out}\| > \epsilon) \leq 2M_H(N)\exp(-2\epsilon^2 N)$$

2.4 Break point

我们希望假设空间 H 的增长函数越小越好，或者至少不要增长的太快，随着 m 的增大，一定会出现一个 m 使假设空间无法shatter。这种不满足的情况说明增长函数从这个点开始变缓了，所以我们把第一个不满足shatter的 m 值称为break point。其准确定义如下：例如，在二维平面中，有1个点时2条

if no k inputs can be shattered by \mathcal{H} ,
call k a **break point** for \mathcal{H}

- $m_{\mathcal{H}}(k) < 2^k$
- $k + 1, k + 2, k + 3, \dots$ also break points!
- will study **minimum break point** k

图 4: Break Point的定义

线，2个点4条线，3个点8条线，但是四个点就只有14条线而不是16条线，因此4就是二维平面的 break point. Break Point 表明增长函数在 k 这个点开始变缓,将 $M_H(N)$ 的上界变小或者有界，使得学习可行¹。

3 总结

VC 维涉及的内容和概念比较广泛，本人只是介绍了Hoeffding不等式在机器学习中应用，讨论了假设函数趋于无限多的情况，引出解决该问题的一个概念 break point.

4 参考文献

- [1] 《台大机器学习基石》 Hoeffding不等式
- [2] <https://www.zhihu.com/question/38607822/answer/149407083>
- [3] <https://www.jianshu.com/p/97a4d3c35aa8>

¹李雨佳-VC学习笔记