

# Understanding Machine Learning

Yuyang Zhang

2017年7月16日

# 目录

<b>1</b>	<b>why can machine learn</b>	<b>3</b>
1.1	基本符号 . . . . .	3
1.1.1	输入 . . . . .	3
1.1.2	输出 . . . . .	3
1.1.3	数据生成模型 . . . . .	3
1.1.4	衡量标准 . . . . .	3
1.2	经验风险最小化 . . . . .	4
1.2.1	过拟合 . . . . .	4
1.3	归纳偏好 . . . . .	4
1.4	为什么可以学习到东西 . . . . .	5
1.5	思路总结 . . . . .	7
<b>2</b>	<b>Probably Approximately Correct</b>	<b>9</b>
2.1	PAC学习理论 . . . . .	9
2.1.1	PAC可学习 . . . . .	9
2.1.2	采样复杂度 . . . . .	9
2.2	泛化PAC理论 . . . . .	9
2.3	不可知PAC学习 . . . . .	10
2.4	学习问题建模 . . . . .	10
2.4.1	广义损失函数 . . . . .	11
2.5	思路总结 . . . . .	11
<b>3</b>	<b>附录1: 不等式证明</b>	<b>12</b>
3.1	Markov Inequality . . . . .	12

# 1 why can machine learn

## 1.1 基本符号

### 1.1.1 输入

Domain Set: 一个任意集合 $\mathcal{X}$ ，也作领域集。可以理解为所有样本的集合，其中每个样本通常以一个能够表征其特征的向量表示。

Label Set: 标签集 $\mathcal{Y}$ 。样本所属于的类别，通常二分类问题，标签集为 $\{0, 1\}$ 或者是 $\{-1, +1\}$ 。

Training Data: 训练数据 $\mathcal{S} = \{(x_1, y_1), \dots, (x_m, y_m)\}$ ，也叫训练集，样本集。

### 1.1.2 输出

Predicting Rule: 预测规则， $h : \mathcal{X} \rightarrow \mathcal{Y}$ 。该规则是一个由样本集到标签集的映射，可以理解为预测器(predictor)，假设(hypothesis)，分类器(classifier)，等等。

### 1.1.3 数据生成模型

我们假定样本 $\mathcal{S}$ 是由概率分布 $\mathcal{D}$ 生成，并且根据一个标记函数 $f : \mathcal{X} \rightarrow \mathcal{Y}$ ，来标记样本的类别，对于任意 $i = 1, \dots, m$ 都有  $y_i = f(x_i)$ ，然而我们并不知道概率分布 $\mathcal{D}$ 与标记函数 $f$ ，我们的目标就是找一个合适的假设 $h$ ，令其与标记函数 $f$ 可以对样本做出相同的标记。

### 1.1.4 衡量标准

我们定义分类误差为：未能成功预测随机数据点正确标签的概率，即对于随机的一个 $x \in \mathcal{X}$ ， $h(x) \neq f(x)$ 的概率。

定义 $\mathcal{A} \subseteq \mathcal{X}$ 为一个领域子集， $\mathcal{A}$ 中的任意实例 $x \in \mathcal{A}$ 的出现概率由 $\mathcal{D}(\mathcal{A})$ 所决定。通常，我们称 $\mathcal{A}$ 为一个事件， $\mathcal{A} = \{x \in \mathcal{X} : \pi(x) = 1\}$ ，其中 $\pi : \mathcal{X} \rightarrow \{0, 1\}$ ，表示样本是否被观测到。我们也将 $\mathcal{D}(\mathcal{A})$ 写作 $\mathbb{P}_{x \sim \mathcal{D}}[\pi(x)]$ 。此时我们可以定义假设 $h$ 的错误率为：

$$\mathcal{L}_{\mathcal{D}, f}(h) \stackrel{\text{def}}{=} \mathbb{P}_{x \sim \mathcal{D}}[h(x) \neq f(x)] \stackrel{\text{def}}{=} \mathcal{D}(\{x | h(x) \neq f(x)\})$$

其中误差的测量是基于概率分布 $\mathcal{D}$ 和标记函数 $f$ 的， $\mathcal{L}_{\mathcal{D},f}(h)$ 也称为泛化误差，**损失**或者 $h$ 的真实误差。

## 1.2 经验风险最小化

机器学习的过程都是基于训练集 $\mathcal{S}$ 的，训练集 $\mathcal{S}$ 由未知分布 $\mathcal{D}$ 从领域集 $\mathcal{X}$ 中采样得出，并由标记函数 $f$ 标记，机器学习的输出是一个基于训练集 $\mathcal{S}$ 的假设， $h_{\mathcal{S}} : \mathcal{X} \rightarrow \mathcal{Y}$ 。

由于我们并不知道分布 $\mathcal{D}$ 与标记函数 $f$ ，所以我们只能根据训练集 $\mathcal{S}$ 来判断我们所选择假设的表现，定义训练误差为：

$$\mathcal{L}_{\mathcal{S}} \stackrel{\text{def}}{=} \frac{|\{x_i | h(x_i) \neq y_i, i = 1, \dots, m\}|}{m}$$

训练误差也称作经验误差和经验风险。

因为训练集 $\mathcal{S}$ 是领域 $\mathcal{X}$ 的一个子集，所以训练样本是真实世界的一个缩影，正如概率统计的一个核心思想，通过样本反映总体。所以我们认为利用样本集寻找一个较好的假设是可行的，即最小化训练误差 $\mathcal{L}_{\mathcal{S}}(h)$ ，这称之为经验风险最小化，ERM(Experience Risk Minimize)。

### 1.2.1 过拟合

一个假设在训练集上效果优异，但在真实世界中表现却糟糕，这种现象称之为过拟合。当我们过度追求经验风险最小化原则时，我们就有可能面临过拟合的风险。

## 1.3 归纳偏好

虽然经验风险最小化会有过拟合的风险，相比于抛弃这个原则，我们更愿意去修正这个原则，考虑我们对于假设的归纳偏好。我们通常的解决方案是根据我们定好的归纳偏好，在有限的假设空间中去搜索所要用的假设，这些假设的集合成为假设类，记为 $\mathcal{H}$ ，则我们的学习过程可以记为：

$$ERM_{\mathcal{H}}(\mathcal{S}) \in \arg \max_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{S}}(h)$$

我们常见的正则化，就是一种归纳偏好， $L1$ 正则化表明我们的归纳偏好是更喜欢参数稀疏的假设。

### 1.4 为什么可以学习到东西

本节旨在说明：当拥有足够多样本时，在有限假设空间 $\mathcal{H}$ 中，经验风险最小化 $ERM_{\mathcal{H}}$ 原则不会出现过拟合，即我们通过训练样本，可以找到一个足够好的假设 $h_s$ ，在真实世界中表现也足够好。

**定义1.1 可实现性假设** 存在 $h^* \in \mathcal{H}$ ，使得 $\mathcal{L}_{\mathcal{D},f}(h^*) = 0$ 。

该假设意味着，对于随机样本集 $\mathcal{S}$ ，由概率分布 $\mathcal{D}$ 采样，由标记函数 $f$ 标记，以概率1使得 $L_{\mathcal{S}}(h^*) = 0$ ，其中样本集 $\mathcal{S}$ 中的样本是独立同分布的。我们定义 $h_{\mathcal{S}}$ 为对 $\mathcal{S}$ 利用 $ERM_{\mathcal{H}}$ 得到的结果：

$$h_{\mathcal{S}} \in \arg \max_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{S}}(h)$$

因为样本集 $\mathcal{S}$ 仍是领域集 $\mathcal{X}$ 的子集，是根据分布随机得到的实例集合，会有一定概率使得采样得到的样本不具有代表性，并不能反映真实的总体情况，所以根据样本集 $\mathcal{S}$ 得到的假设 $h_{\mathcal{S}}$ 并不一定准确，在真实世界中的表现有可能很差。因此，我们选择一定程度的容忍，容忍会有一定几率采样到不具有代表性的样本，一般来说，我们定义采样得到不具有代表性样本的概率之多为 $\delta$ ，则 $1 - \delta$ 为置信参数。同时对于假设的预测效果，我们能容忍的损失上限为 $\epsilon$ ，称为精度参数，如果 $\mathcal{L}_{\mathcal{D},f}(h_{\mathcal{S}}) > \epsilon$ ，那么这是一个差的假设，反之，则是一个好的假设。

对于我们的样本集 $\mathcal{S}$ ，最好的假设 $h_{\mathcal{S}}$ 仍有 $\mathcal{L}_{\mathcal{D},f}(h_{\mathcal{S}}) > \epsilon$ ，则这组样本采样失败，因为我们的 $ERM_{\mathcal{H}}$ 无法从 $\mathcal{S}$ 中学到有用的东西。当样本集中所有样本 $(x_1, \dots, x_m)$ ，都不具有代表性，我们记为：

$$\{x_i | x_i \in \mathcal{S}, i = 1, \dots, m, \mathcal{L}_{\mathcal{D},f}(h_{\mathcal{S}}) > \epsilon\}$$

则采样失败的概率上界为：

$$\mathcal{D}^m(\{x_i | x_i \in \mathcal{S}, i = 1, \dots, m, \mathcal{L}_{\mathcal{D},f}(h_{\mathcal{S}}) > \epsilon\})$$

因为每个样本都是不具有代表性的样本，所以是采样失败的概率的上界。设 $\mathcal{H}_B$ 为差的假设的集合：

$$\mathcal{H}_B = \{h | \mathcal{L}_{\mathcal{D},f}(h) > \epsilon, h \in \mathcal{H}\}$$

同时设：

$$M = \{x | x \in \mathcal{S}, \exists h \in \mathcal{H}_B, \mathcal{L}_{\mathcal{S}}(h) = 0\}$$

为误导集，误导集使差的假设在训练样本 $\mathcal{S}$ 上表现良好，但在真实世界中表现较差。因为有可实现假设

$$\exists h^*, \mathcal{L}_{\mathcal{D},f}(h^*) = 0$$

且有

$$h_{\mathcal{S}} \in \arg \max_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{S}}(h)$$

产生 $\mathcal{L}_{\mathcal{D},f}(h_{\mathcal{S}}) > \epsilon$  是因为样本集不好，样本集采样失败，所以，当且仅当 $\mathcal{S} \subseteq M$ 时，才会出现 $\mathcal{L}_{\mathcal{D},f}(h_{\mathcal{S}}) > \epsilon$ ，我们将其表示为：

$$\{x_i | x_i \in \mathcal{S}, i = 1, \dots, m, \mathcal{L}_{\mathcal{D},f}(h_{\mathcal{S}}) > \epsilon\} \subseteq M$$

所以我们有

$$\mathcal{D}^m(\{x_i | x_i \in \mathcal{S}, i = 1, \dots, m, \mathcal{L}_{\mathcal{D},f}(h_{\mathcal{S}}) > \epsilon\}) \leq \mathcal{D}^m(M)$$

而 $M$ 又可以写作：

$$M = \bigcup_{h \in \mathcal{H}_B} \{x | x \in \mathcal{S}, \mathcal{L}_{\mathcal{S}}(h) = 0\}$$

则：

$$\mathcal{D}^m(\{x_i | x_i \in \mathcal{S}, i = 1, \dots, m, \mathcal{L}_{\mathcal{D},f}(h_{\mathcal{S}}) > \epsilon\}) \leq \mathcal{D}^m(\bigcup_{h \in \mathcal{H}_B} \{x | x \in \mathcal{S}, \mathcal{L}_{\mathcal{S}}(h) = 0\})$$

由 $P(A \cup B) \leq P(A) + P(B)$ 得：

$$\mathcal{D}^m(\bigcup_{h \in \mathcal{H}_B} \{x | x \in \mathcal{S}, \mathcal{L}_{\mathcal{S}}(h) = 0\}) \leq \sum_{h \in \mathcal{H}_B} \mathcal{D}^m(\{x | x \in \mathcal{S}, \mathcal{L}_{\mathcal{S}}(h) = 0\})$$

其中我们将 $\mathcal{S}$ 拆开，

$$\mathcal{D}^m(\{x | x \in \mathcal{S}, \mathcal{L}_{\mathcal{S}}(h) = 0\}) = \mathcal{D}^m(\{x_i | h(x_i) = f(x_i), x_i \in \mathcal{S}, i = 1, \dots, m\})$$

$$\mathcal{D}^m(\{x | x \in \mathcal{S}, \mathcal{L}_{\mathcal{S}}(h) = 0\}) = \prod_{i=1}^m \mathcal{D}(\{x_i | h(x_i) = f(x_i), x_i \in \mathcal{S}\})$$

根据 $1 - \epsilon \leq e^{-\epsilon}$ ，对于等号右边的连乘的每一项都有：

$$\mathcal{D}(\{x_i | h(x_i) = y_i, x_i \in \mathcal{S}\}) = 1 - \mathcal{L}_{\mathcal{D},f}(h) \leq 1 - \epsilon$$

对于所有的  $h \in \mathcal{H}_B$ :

$$\mathcal{D}^m(\{x|x \in \mathcal{S}, \mathcal{L}_{\mathcal{S}}(h) = 0\}) \leq (1 - \epsilon)^m \leq e^{-\epsilon m}$$

可得:

$$\mathcal{D}^m(\{x_i|x_i \in \mathcal{S}, i = 1, \dots, m, \mathcal{L}_{\mathcal{D},f}(h_{\mathcal{S}}) > \epsilon\}) \leq |\mathcal{H}_B|e^{-\epsilon m} \leq |\mathcal{H}|e^{-\epsilon m}$$

所以, 若我们对采样失败概率的容忍大于采样失败的概率上限, 我们认为最后得到的  $h_{\mathcal{S}}$  学到了有用的东西的, 记为:

$$|\mathcal{H}|e^{-\epsilon m} \leq \delta$$

借此我们可以推断出我们需要拥有的样本数量为:

$$m \geq \frac{\ln(|\mathcal{H}|/\delta)}{\epsilon}$$

**推论1.1** 设  $\mathcal{H}$  为一个有限假设集合,  $\delta \in (0, 1)$ ,  $\epsilon > 0$ , 当

$$m \geq \frac{\ln(|\mathcal{H}|/\delta)}{\epsilon}$$

成立时, 从而对于任何标记函数  $f$ 、任何分布  $\mathcal{D}$ , 可实现性假设最少以  $1 - \delta$  的概率, 对于每个  $ERM$  假设  $h_{\mathcal{S}}$ , 有以下不等式成立:

$$\mathcal{L}_{\mathcal{D},f}(h_{\mathcal{S}}) \leq \epsilon$$

上述推论表明, 对于足够大的  $m$ , 由  $ERM_{\mathcal{H}}$  规则生成的有限假设会以  $1 - \delta$  概率得到小于误差  $\epsilon$  的近似正确解, 概率在于容忍采样失败概率  $\delta$ , 近似在于容忍误差  $\epsilon$ , 即我们的算法, 可能 (容忍了失败概率) 会学到令我们基本满意 (容忍了一定误差) 的东西。而概率近似正确, 即 PAC 理论, 将在第二章详述。

## 1.5 思路总结

1. 为什么不能有效学习? 为什么会  $\mathcal{L}_{\mathcal{D},f}(h) > \epsilon$ ?
2. 因为样本不具有代表性, 学到的东西有问题, 采样失败。
3. 怎么解决?

4. 降低采样失败概率至我们可以容忍的范围，小于 $\delta$ 。
5. 采样失败的概率上限小于 $\delta$ 。
6.  $\sup \mathcal{D}^m(\{x_i | x_i \in \mathcal{S}, i = 1, \dots, m, \mathcal{L}_{\mathcal{D}, f}(h_{\mathcal{S}}) > \epsilon\}) = |\mathcal{H}|e^{-\epsilon m} \leq \delta$
7.  $m \geq \frac{\ln(|\mathcal{H}|/\delta)}{\epsilon}$
8. 当样本足够充足时，学习到的假设概率近似正确。



## 2 Probably Approximately Correct

### 2.1 PAC学习理论

#### 2.1.1 PAC可学习

接着上一章继续说，在经验风险最小化准则下，对于一个有限假设类，如果有足够多的训练样本，则我们输出的学习算法在真实世界中，是概率近似正确的。我们做以下定义：

**定义2.1 PAC可学习** 若存在一个函数  $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathcal{N}$  和一个学习算法，使得对于给定的  $\epsilon, \delta \in (0, 1)$  和任一分布  $\mathcal{D}$ 、任一标记函数  $f : \mathcal{X} \rightarrow \{0, 1\}$ ，可实现假设成立时，那么当样本数量  $m \geq m_{\mathcal{H}}(\epsilon, \delta)$  时，算法将以不小于  $1 - \delta$  的概率返回一个使

$$\mathcal{L}_{\mathcal{D}, f}(h) \leq \epsilon$$

的假设  $h$ 。

#### 2.1.2 采样复杂度

函数  $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathcal{N}$  决定了假设的采样复杂度，即可以学到东西时所需要的样本数量。采样复杂度不仅依赖于  $\epsilon$  和  $\delta$ ，同样还依赖于假设空间中的假设数量：

$$m_{\mathcal{H}} = \frac{\ln(|\mathcal{H}|/\delta)}{\epsilon}$$

其与假设数量的对数成正比。

**推论2.1** 当采样复杂度满足：

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{\ln(|\mathcal{H}|/\delta)}{\epsilon} \right\rceil$$

则任一有限类假设是PAC可学习的。

$m_{\mathcal{H}}(\epsilon, \delta)$  被称为假设结合  $\mathcal{H}$  的采样复杂度的最小函数。由上述定义我们可知，一个假设是否是PAC可学习的，还依赖于假设空间的大小，而假设空间则与VC维相关，将会在后几章说明。

### 2.2 泛化PAC理论

为了让该理论更贴近实际，则我们考虑将理论的前提约束放宽，进行泛化。

1. 去掉可实现性假设。
2. 考虑多分类、回归等问题。

### 2.3 不可知PAC学习

当我们放弃可实现假设时，我们称之为不可知，具体可以理解为：完全一样的样本，却又不同的标记，这种情况下，我们的算法如何学到东西。

从此处开始，为了简写，我们将  $\mathcal{D}$  定义为  $\mathcal{X} \times \mathcal{Y}$  上的概率分布，即将之前  $\mathcal{D}$  和  $f$  简写在一起，作为领域集和标签集的联合概率分布，定义真实误差为：

$$\mathcal{L}_{\mathcal{D}}(h) \stackrel{\text{def}}{=} \mathbb{P}_{(x,y) \sim \mathcal{D}}[h(x) \neq y] \stackrel{\text{def}}{=} \mathcal{D}(\{(x,y) : h(x) \neq y\})$$

在面对不可知的情况下，最好的预测器为贝叶斯预测器：

$$f_{\mathcal{D}}(x) = \begin{cases} 1 & \mathbb{P}[y = 1|x] \geq \frac{1}{2} \\ 0 & \text{else} \end{cases}$$

我们易证贝叶斯预测器是最优的：对于每个领域集的实例  $x \in \mathcal{X}$  都选择数量较多的类别，即每个实例的分类正确概率都大于  $1/2$ ，没有其他分类器可以达到比它更低的错误率，对于任意预测器  $g$  都有  $\mathcal{L}_{\mathcal{D}}(f_{\mathcal{D}}) \leq \mathcal{L}_{\mathcal{D}}(g)$ ，所以该预测器最优。但我们并不知道实际的概率分布  $\mathcal{D}$ ，所以我们并没法使用这样的预测器。

所以对于不可知问题，我们只能选择容忍（看来统计的机器学习就是一个容忍的过程Orz）：

**定义2.2 不可知PAC可学习** 若存在一个函数  $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathcal{N}$  和一个学习算法，使得对于给定的  $\epsilon, \delta \in (0, 1)$  和任一分布  $\mathcal{D}$ ，当样本数量  $m \geq m_{\mathcal{H}}(\epsilon, \delta)$  时，算法将以不小于  $1 - \delta$  的概率返回一个使

$$\mathcal{L}_{\mathcal{D}}(h) \leq \min_{h' \in \mathcal{H}} \mathcal{L}_{\mathcal{D}}(h') + \epsilon$$

成立的假设  $h$ 。

### 2.4 学习问题建模

此节我们讨论学习问题建模，总的来说，除了二分类问题，学习任务还分为以下几种：

- 多分类
- 回归

虽然说学习任务分为多种，但主要区别只在于损失函数不同。

### 2.4.1 广义损失函数

给定任意集合 $\mathcal{H}$ 和定义域 $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ ，令 $\ell$ 为 $\mathcal{H} \times \mathcal{Z}$ 到非负实数的映射，记为 $\ell: \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$ ， $\ell$ 就是损失函数。现在我们重新定义真实误差和经验误差：

$$\begin{aligned}\mathcal{L}_{\mathcal{D}}(h) &\stackrel{\text{def}}{=} \mathbb{E}_{z \sim \mathcal{D}}[\ell(h, z)] \\ \mathcal{L}_{\mathcal{S}}(h) &\stackrel{\text{def}}{=} \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)\end{aligned}$$

同时我们定义常用的损失函数：

- 0-1损失

$$\ell_{0-1}(h, (x, y)) \stackrel{\text{def}}{=} \begin{cases} 0 & h(x) = y \\ 1 & h(x) \neq y \end{cases}$$

- 平方损失

$$\ell_{sq}(h, (x, y)) \stackrel{\text{def}}{=} (h(x) - y)^2$$

**定义2.3 广义损失函数下的不可知PAC可学习** 对于集合 $\mathcal{Z}$ 和损失函数 $\ell: \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$ ，若存在一个函数 $m_{\mathcal{H}}: (0, 1)^2 \rightarrow \mathbb{N}$ 和一个学习算法，使得对于给定的 $\epsilon, \delta \in (0, 1)$ 和任一分布 $\mathcal{D}$ ，当样本数量 $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ 时，算法将以不小于 $1 - \delta$ 的概率返回一个使

$$\mathcal{L}_{\mathcal{D}}(h) \leq \min_{h' \in \mathcal{H}} \mathcal{L}_{\mathcal{D}}(h') + \epsilon$$

成立的假设 $h$ ，其中 $\mathcal{L}_{\mathcal{D}}(h) = \mathbb{E}_{z \sim \mathcal{D}}[\ell(h, z)]$ 。

## 2.5 思路总结

本章主要在于泛化PAC学习的试用范围，通过更改假设来让这个理论更具有普适性，更贴近生活现实。（通俗来说就是通过不断的容忍，放缩范围，看来搞统计的都是受啊XD）

### 3 附录1: 不等式证明

#### 3.1 Markov Inequality

马尔科夫不等式:

$$\mathbb{P}(X \geq a) \leq \frac{\mathbb{E}(X)}{a}$$

其中 $X$ 是非负随机变量。

证:

$$\begin{aligned} \mathbb{E}(X) &= \int_0^{+\infty} x f(x) dx \\ &= \int_0^a x f(x) dx + \int_a^{+\infty} x f(x) dx \\ &\leq \int_0^a 0 f(x) dx + \int_a^{+\infty} a f(x) dx \\ &\leq 0 + \int_a^{+\infty} a f(x) dx = a \int_a^{+\infty} f(x) dx = a \mathbb{P}(X > a) \end{aligned}$$

移动一下 $a$ 的位置, 不等式得证, 其中第二行到第三行是将积分中的 $x$ 换成积分下限。