

# 机器学习（西瓜书） 注 解

（第 8 章 集成学习）

<https://blog.csdn.net/jbb0523>

# 前言

经常听人说南大周老师所著的《机器学习》（以下统称为西瓜书）是一本入门教材，是一本科普性质的教科书。在该书第十次印刷之际，周老师在“[如何使用本书](#)”中也提到“这是一本入门级教科书”。然而，本人读起来却感觉该书远不止“科普”“入门”那么简单，书中的很多公式需要思考良久方能推导，很多概念需要反复咀嚼才能消化。边读边想着是不是应该将自己学习时遇到的一些知识难点的解析分享出来，以帮助更多的人入门。自己的确也随手做过一些笔记，但由于怀疑这仅是自己的个别现象，毕竟读书期间，思考更多的是如何使用单片机、DSP、ARM、FPGA 等，而这些基本是不需要推导任何公式的，因此作罢。偶然间在[周老师的新浪微博](#)看到如下对话：



此时方知，可能“读不懂”并不是个别现象。因此决定写一本“西瓜书注解”或者称为“西瓜书读书笔记”，对自己研读西瓜书时遇到的“台阶”进行解释和推导，以帮助更多的人能够更快地进入到这个领域。另外，近期越来越强地意识到，扎扎实实地推导一些基础算法的公式，无论是对于理解算法本身机理还是进行学术研究，都是非常有必要的。

自己会根据个人学习进度和研究需要按章发布，不知道能不能坚持写完，加油！

毕竟自己也是一名初学者，所以可能一些概念解释并不完整、一些公式推导并不优美，甚至会存在错误，这是不可避免的，不接受谩骂，但欢迎将问题反馈给我，共同学习进步！

（网盘链接：<https://pan.baidu.com/s/1QtEiNnk8jMzmbs0KPBN-w>）

# 第 8 章目录

第 8 章 集成学习.....	1
8.1 个体与集成.....	1
1、式(8.1)的解释 .....	1
2、式(8.2)的解释 .....	2
3、式(8.3)的推导 .....	2
8.2 Boosting .....	2
1、式(8.4)的解释 .....	3
2、式(8.5)的解释 .....	3
3、式(8.6)的推导 .....	3
4、式(8.7)的推导 .....	4
5、式(8.8)的推导 .....	5
6、式(8.9)的推导 .....	5
7、式(8.12)的推导[?].....	6
8、式(8.13)的推导 .....	6
9、式(8.14)的推导[?].....	6
10、式(8.16)的推导 .....	7
11、式(8.17)的推导 .....	7
12、式(8.18)的推导 .....	7
13、式(8.19)的推导 .....	8
14、AdaBoost 的个人推导.....	8
15、更深一步理解权重更新公式的含义.....	11
16、能够接受带权样本的基学习算法.....	12
8.3 Bagging 与随机森林 .....	13
1、Bagging 算法的解释 .....	13
2、式(8.20)的解释 .....	14
3、式(8.21)的推导 .....	14
4、随机森林的解释.....	14
8.4 结合策略.....	15
1、硬投票和软投票的解释.....	15
2、元学习器(meta-learner)的解释.....	15
3、Stacking 算法的解释.....	15
8.5 多样性.....	15
1、式(8.28)的解释 .....	15
2、式(8.31)与式(8.28)的等价性推导.....	16
3、式(8.32)的解释 .....	17
4、列联表(contingency table)的解释 .....	17
5、式(8.40)的解释 .....	17
6、式(8.41)的解释 .....	17
7、式(8.42)的解释 .....	17
8、多样性增强的解释.....	18
8.6、GB(Gradient Boosting)/GBDT/XGBoost .....	18

1、梯度下降法.....	18
2、从梯度下降的角度解释 AdaBoost.....	20
3、梯度提升(Gradient Boosting).....	22
4、梯度提升树(GBDT).....	23
5、XGBoost.....	23
8.7 本章小节.....	24

## 第 8 章 集成学习

首先，我们来看一下西瓜书作者的[谷歌学术主页](#)（2018-11-18）：

标题	引用次数	年份
<a href="#">Top 10 algorithms in data mining</a> X Wu, V Kumar, JR Quinlan, J Ghosh, Q Yang, H Motoda, GJ McLachlan, ... Knowledge and information systems 14 (1), 1-37	4149	2008
<a href="#">Ensembling neural networks: many could be better than all</a> ZH Zhou, J Wu, W Tang Artificial intelligence 137 (1-2), 239-263	1773	2002
<a href="#">ML-KNN: A lazy learning approach to multi-label learning</a> ML Zhang, ZH Zhou Pattern recognition 40 (7), 2038-2048	1614	2007
<a href="#">Exploratory undersampling for class-imbalance learning</a> XY Liu, J Wu, ZH Zhou IEEE Trans. Systems, Man, and Cybernetics, Part B: Cybernetics 39 (2), 539-550	995	2009
<a href="#">Ensemble Methods: Foundations and Algorithms</a> ZH Zhou Chapman & Hall/CRC Press	939	2012
<a href="#">A Review on Multi-Label Learning Algorithms</a> ML Zhang, ZH Zhou IEEE Trans. Knowledge and Data Engineering 26 (8), 1819-1837	926	2014
<a href="#">Automatic age estimation based on facial aging patterns</a> X Geng, ZH Zhou, K Smith-Miles IEEE Trans. Pattern Analysis and Machine Intelligence 29 (12), 2234-2240	835	2007

引用次数前 5 名中，除第 3 名“ML-KNN: A lazy learning approach to multi-label learning”与集成学习(ensemble learning)无关以外，其余均与集成学习有直接关系。其中：

第 1 名的“Top 10 algorithms in data mining”是在 ICDM'06 中投票选出的数据挖掘十大算法，每个提名算法均由业内专家代表去阐述，然后进行投票，其中最终得票排名第 7 位的“Adaboost”即由西瓜书作者作为代表进行阐述；第 2 名的“Ensembling neural networks: many could be better than all”是西瓜书作者最早的代表作，该文催生了基于优化的集成修剪(Ensemble pruning)技术；第 4 名的“Exploratory undersampling for class-imbalance learning”是以集成学习技术解决类别不平衡问题；第 5 名的“Ensemble Methods: Foundations and Algorithms”则是西瓜书作者所著的集成学习专著。

因此，虽然周老师现在的研究方向已涉及到了机器学习的方方面面，但若将集成学习看作是其学术之路的源头，实际也不为过。

所以，尽情享受西瓜书作者最拿手的本章内容吧^\_^

### 8.1 个体与集成

本节内容通俗易懂，几乎不需要什么注解。

基学习器(base learner)的概念在论文中经常出现，可留意一下；另外，本节提到的投票法有两种，除了本节的多数投票(majority voting)，还有概率投票(probability voting)，这两点在 8.4 节中均会提及，即硬投票和软投票。

#### 1、式(8.1)的解释

注意到 $f(\mathbf{x})$ 为 $\mathbf{x}$ 的真实标记，而 $h_i(\mathbf{x})$ 为基分类器 $h_i$ 对 $\mathbf{x}$ 的预测标记，因此 $h_i(\mathbf{x}) \neq f(\mathbf{x})$

表示预测标记与真实标记不相同的情况，即预测出错， $P(h_i(\mathbf{x}) \neq f(\mathbf{x}))$ 即为错误率。

## 2、式(8.2)的解释

注意到当前仅针对二分类问题 $y \in \{-1, +1\}$ ，即预测标记 $h_i(\mathbf{x}) \in \{-1, +1\}$ 。因此，若 $T$ 个基分类器中预测输出+1占半数以上，则 $\sum_{i=1}^T h_i(\mathbf{x}) > 0$ ；反之，若 $T$ 个基分类器中预测输出-1占半数以上，则 $\sum_{i=1}^T h_i(\mathbf{x}) < 0$ 。函数 $\text{sign}(\cdot)$ 为符号函数，在 $\cdot < 0$ 时取值为-1，在 $\cdot > 0$ 时取值为+1，即本式的函数 $H(\mathbf{x})$ 与多数基分类器输出一致，也就是多数投票。

## 3、式(8.3)的推导

本式推导基于 Hoeffding 不等式。首先将 P192 习题 8.1 的式(8.43)变形如下：

$$P(H(T) \leq \lfloor T/2 \rfloor) = \sum_{k=0}^{\lfloor T/2 \rfloor} \binom{T}{k} (1-\epsilon)^k \epsilon^{T-k}$$

从式(8.43)到上式的符号变化包括 $n \rightarrow T, k \rightarrow \lfloor T/2 \rfloor, i \rightarrow k, p \rightarrow (1-\epsilon)$ 。相应的，对于接下来的式(8.44)，则对 $\delta > 0, \lfloor T/2 \rfloor = (1-\epsilon-\delta)T$ ，有 Hoeffding 不等式

$$P(H(T) \leq (1-\epsilon-\delta)T) = e^{-2\delta^2 T}$$

对于 $\lfloor T/2 \rfloor = (1-\epsilon-\delta)T$ ，整理变形得 $\delta$ 表达式

$$\delta = 1 - \epsilon - \frac{\lfloor T/2 \rfloor}{T}$$

注意到集成学习要求基分类器泛化性能优于随机猜测的学习器，例如在二分类问题上精度略高于 50% 的分类器，即 $\epsilon < 0.5$ ，因此 $1-\epsilon > \frac{1}{2}$ ；而 $\frac{\lfloor T/2 \rfloor}{T} \leq \frac{1}{2}$ （符号 $\lfloor \cdot \rfloor$ 表示向下取整，“=”在 $T$ 为偶数时成立，若 $T$ 为奇数，则一定是“<”）。因此

$$\delta = 1 - \epsilon - \frac{\lfloor T/2 \rfloor}{T} \geq 1 - \epsilon - \frac{1}{2} > 0$$

函数 $f(x) = e^{-Cx^2}$ 为单调递减函数，其中 $C > 0$ 为常数，因此

$$P(H(T) \leq \lfloor T/2 \rfloor) = e^{-2\delta^2 T} \leq e^{-2(1-\epsilon-\frac{1}{2})^2 T} = e^{-\frac{T}{2}(1-2\epsilon)^2}$$

而根据式(8.2)中 $H(\mathbf{x})$ 的含义， $H(T) \leq \lfloor T/2 \rfloor$ （即 $T$ 个基分类器中预测正确的个数不超过 $\lfloor T/2 \rfloor$ ，也就是预测错误）与 $H(\mathbf{x}) \neq f(\mathbf{x})$ 等价，推导完毕！

注意，上式推导中要保证 $\delta > 0$ ，如此才能保证 $\delta$ 与 $\delta^2$ 符号一致，不等式缩放也一致；习题 8.1 中介绍的 Hoeffding 不等式仅为伯努力随机变量特例，参见 CSDN 博客“[机器学习数学原理（8）——霍夫丁不等式](https://blog.csdn.net/z_x_1996/article/details/73564926)”（[https://blog.csdn.net/z\\_x\\_1996/article/details/73564926](https://blog.csdn.net/z_x_1996/article/details/73564926)）。

## 8.2 Boosting

注意 8.1 节最后一段提到：根据个体学习器的生成方式，目前的集成学习方法大致可分为两大类，即个体学习器间存在强依赖关系、必须串行生成的序列化方法，以及个体学习器间不存在强依赖关系、可同时生成的并行化方法。

本节 Boosting 为前者的代表，Adaboost 又是 Boosting 族算法的代表。

## 1、式(8.4)的解释

注意到 $h_t(\mathbf{x}) \in \{-1, +1\}$ ，加权系数 $\alpha_t$ 由后面的式(8.11)给出（ $T$ 个系数 $\alpha_t$ 仅与自身对应的基学习器错误率 $\epsilon_t$ 有关，相互之间无必然联系；由于要求 $\epsilon_t < 0.5$ ，因此 $\alpha_t > 0$ ），所得线性加权结果 $H(\mathbf{x})$ 就是一个实数，可能大于零，也可能小于零，取值范围无法确定。

## 2、式(8.5)的解释

(1)先考虑指数损失函数 $e^{-f(\mathbf{x})H(\mathbf{x})}$ 的含义（参见 P131 图 6.5）： $f$ 为真实函数，对于样本 $\mathbf{x}$ 来说， $f(\mathbf{x}) \in \{-1, +1\}$ 只能取 $-1$ 和 $+1$ 两个值，而 $H(\mathbf{x})$ 是一个实数；

当 $H(\mathbf{x})$ 的符号与 $f(\mathbf{x})$ 一致时， $f(\mathbf{x})H(\mathbf{x}) > 0$ ，因此 $e^{-f(\mathbf{x})H(\mathbf{x})} = e^{-|H(\mathbf{x})|} < 1$ ，且 $|H(\mathbf{x})|$ 越大指数损失函数 $e^{-f(\mathbf{x})H(\mathbf{x})}$ 越小（这很合理：此时 $|H(\mathbf{x})|$ 越大意味着分类器本身对预测结果的信心越大，损失应该越小；若 $|H(\mathbf{x})|$ 在零附近，虽然预测正确，但表示分类器本身对预测结果信心很小，损失应该较大）；

当 $H(\mathbf{x})$ 的符号与 $f(\mathbf{x})$ 不一致时， $f(\mathbf{x})H(\mathbf{x}) < 0$ ，因此 $e^{-f(\mathbf{x})H(\mathbf{x})} = e^{|H(\mathbf{x})|} > 1$ ，且 $|H(\mathbf{x})|$ 越大指数损失函数 $e^{-f(\mathbf{x})H(\mathbf{x})}$ 越大（这很合理：此时 $|H(\mathbf{x})|$ 越大意味着分类器本身对预测结果的信心越大，但预测结果是错的，因此损失应该越大；若 $|H(\mathbf{x})|$ 在零附近，虽然预测错误，但表示分类器本身对预测结果信心很小，虽然错了，损失应该较小）；

对比 0/1 损失函数（参见 P131 图 6.5）：当 $H(\mathbf{x})$ 的符号与 $f(\mathbf{x})$ 一致时则损失为 0，当 $H(\mathbf{x})$ 的符号与 $f(\mathbf{x})$ 不一致时则损失为 1，而不管 $H(\mathbf{x})$ 的大小。

(2)再解释符号 $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\cdot]$ 的含义： $\mathcal{D}$ 为概率分布，可简单理解为在数据集 $D$ 中进行一次随机抽样，每个样本被取到的概率； $\mathbb{E}[\cdot]$ 为经典的期望，则综合起来 $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\cdot]$ 表示在概率分布 $\mathcal{D}$ 上的期望，可简单理解为对数据集 $D$ 以概率 $\mathcal{D}$ 进行加权后的期望。

综上所述，若数据集 $D$ 中样本 $\mathbf{x}$ 的权值分布为 $\mathcal{D}(\mathbf{x})$ ，则式(8.5)可写为：

$$\begin{aligned}\ell_{\text{exp}}(H \mid \mathcal{D}) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H(\mathbf{x})}] \\ &= \sum_{\mathbf{x} \in D} \mathcal{D}(\mathbf{x}) e^{-f(\mathbf{x})H(\mathbf{x})} \\ &= \sum_{\mathbf{x} \in D} \mathcal{D}(\mathbf{x}) (e^{-H(\mathbf{x})} \mathbb{I}(f(\mathbf{x}) = 1) + e^{H(\mathbf{x})} \mathbb{I}(f(\mathbf{x}) = -1))\end{aligned}$$

特别地，若针对任意样本 $\mathbf{x}$ ，若分布 $\mathcal{D}(\mathbf{x}) = \frac{1}{|D|}$ ，其中 $|D|$ 为数据集 $D$ 样本个数，则

$$\ell_{\text{exp}}(H \mid \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H(\mathbf{x})}] = \frac{1}{|D|} \sum_{\mathbf{x} \in D} e^{-f(\mathbf{x})H(\mathbf{x})}$$

而这就是在求传统平均值。

## 3、式(8.6)的推导

首先，当 $f(\mathbf{x}) = 1$ 时 $e^{-f(\mathbf{x})H(\mathbf{x})} = e^{-H(\mathbf{x})}$ ，当 $f(\mathbf{x}) = -1$ 时 $e^{-f(\mathbf{x})H(\mathbf{x})} = e^{H(\mathbf{x})}$ ；

其次，由以上分析可知式(8.5)包含 $|D|$ 个不同的 $e^{-f(\mathbf{x})H(\mathbf{x})}$ ，其中 $\mathbf{x} \in D$ ；

最后，针对每个 $H(\mathbf{x})$ 求导，只有对应的 $e^{-f(\mathbf{x})H(\mathbf{x})}$ 会起作用，其余与 $H(\mathbf{x})$ 无关的项求

导后等于零，若不理解本句话，将式(8.5)的解释中的式(8.5)改写为如下形式

$$\begin{aligned}\ell_{\text{exp}}(H \mid \mathcal{D}) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H(\mathbf{x})}] = \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i)H(\mathbf{x}_i)} \\ &= \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}(\mathbf{x}_i) (e^{-H(\mathbf{x}_i)} \mathbb{I}(f(\mathbf{x}_i) = 1) + e^{H(\mathbf{x}_i)} \mathbb{I}(f(\mathbf{x}_i) = -1))\end{aligned}$$

但这里有点不同的是，式(8.6)是为了推导出式(8.7)和式(8.8)的判定准则预测 $\mathbf{x}$ 的类别，即无法得知 $f(\mathbf{x})$ 具体等于1还是-1，只能知道 $f(\mathbf{x})$ 为1的后验概率 $P(f(\mathbf{x}) = 1 \mid \mathbf{x})$ 和为-1的后验概率 $P(f(\mathbf{x}) = -1 \mid \mathbf{x})$ （ $f(\mathbf{x})$ 的定义在式(8.1)上方一段话），因此上式应该写为

$$\begin{aligned}\ell_{\text{exp}}(H \mid \mathcal{D}) &= \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}(\mathbf{x}_i) \mathbb{E}[e^{-f(\mathbf{x}_i)H(\mathbf{x}_i)}] \\ &= \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}(\mathbf{x}_i) (e^{-H(\mathbf{x}_i)} P(f(\mathbf{x}_i) = 1 \mid \mathbf{x}_i) + e^{H(\mathbf{x}_i)} P(f(\mathbf{x}_i) = -1 \mid \mathbf{x}_i))\end{aligned}$$

现对 $H(\mathbf{x}_i)$ 求导，求和项中只有第 $i$ 项 $\mathcal{D}(\mathbf{x}_i) \mathbb{E}[e^{-f(\mathbf{x}_i)H(\mathbf{x}_i)}]$ 会起作用。 $\mathcal{D}(\mathbf{x}_i)$ 为常量，不影响接下来式(8.7)中令导数等于零求解的过程，而

$$\frac{\partial e^{-H(\mathbf{x})}}{\partial H(\mathbf{x})} = -e^{-H(\mathbf{x})}, \quad \frac{\partial e^{H(\mathbf{x})}}{\partial H(\mathbf{x})} = e^{H(\mathbf{x})}$$

因此求导结果为

$$\frac{\partial \ell_{\text{exp}}(H \mid \mathcal{D})}{\partial H(\mathbf{x})} = -e^{-H(\mathbf{x})} P(f(\mathbf{x}) = 1 \mid \mathbf{x}) + e^{H(\mathbf{x})} P(f(\mathbf{x}) = -1 \mid \mathbf{x})$$

即式(8.6)。

#### 4、式(8.7)的推导

$$\text{令式(8.6)等于零: } -e^{-H(\mathbf{x})} P(f(\mathbf{x}) = 1 \mid \mathbf{x}) + e^{H(\mathbf{x})} P(f(\mathbf{x}) = -1 \mid \mathbf{x}) = 0$$

$$\text{移项: } e^{H(\mathbf{x})} P(f(\mathbf{x}) = -1 \mid \mathbf{x}) = e^{-H(\mathbf{x})} P(f(\mathbf{x}) = 1 \mid \mathbf{x})$$

$$\text{两边同乘 } \frac{e^{H(\mathbf{x})}}{P(f(\mathbf{x})=-1|\mathbf{x})}: \quad e^{2H(\mathbf{x})} = \frac{P(f(\mathbf{x})=1|\mathbf{x})}{P(f(\mathbf{x})=-1|\mathbf{x})}$$

$$\text{取} \ln(\cdot): \quad 2H(\mathbf{x}) = \ln \frac{P(f(\mathbf{x})=1|\mathbf{x})}{P(f(\mathbf{x})=-1|\mathbf{x})}$$

两边同除 $\frac{1}{2}$ 即得式(8.7)。

式(8.6)和式(8.7)的推导在原始文献[Friedman J H, Hastie T, Tibshirani R, et al. [Additive logistic regression : A statistical view of boosting](#)[J]. Annals of Statistics, 2000, 28(2): 337-407.]

第9页(第345页)也有推导:



LEMMA 1.  $E(e^{-yF(x)})$  is minimized at

$$(12) \quad F(x) = \frac{1}{2} \log \frac{P(y = 1|x)}{P(y = -1|x)}.$$

Hence

$$(13) \quad P(y = 1|x) = \frac{e^{F(x)}}{e^{-F(x)} + e^{F(x)}},$$

$$(14) \quad P(y = -1|x) = \frac{e^{-F(x)}}{e^{-F(x)} + e^{F(x)}}.$$

PROOF. While  $E$  entails expectation over the joint distribution of  $y$  and  $x$ , it is sufficient to minimize the criterion conditional on  $x$ :

$$\begin{aligned} E(e^{-yF(x)}|x) &= P(y = 1|x)e^{-F(x)} + P(y = -1|x)e^{F(x)}, \\ \frac{\partial E(e^{-yF(x)}|x)}{\partial F(x)} &= -P(y = 1|x)e^{-F(x)} + P(y = -1|x)e^{F(x)}. \end{aligned}$$

The result follows by setting the derivative to zero.  $\square$

## 5、式(8.8)的推导

该式通俗易懂，这里仅解释一下贝叶斯错误率的概念。这来源于 P148 的式(7.6)表示的贝叶斯最优分类器，可以发现式(8.8)的最终结果是式(7.6)的二分类特殊形式。

到此为止，本节证明了指数的损失函数是分类任务原本 0/1 损失函数的一致替代损失函数（参见 P131 图 6.5），而指数损失函数有更好的数学性质，例如它是连续可微函数，因此接下来的式(8.9)至式(8.19)基于指数损失函数推导 AdaBoost 的理论细节。

## 6、式(8.9)的推导

【注 1】本节符号略有混乱，如前所述式(8.4)的  $H(\mathbf{x})$  是连续实值函数，但在图 8.3 最后一行的输出  $H(\mathbf{x})$  明显只能取  $-1$  和  $+1$  两个值（与式(8.2)相同），本节除了图 8.3 最后一行的输出之外， $H(\mathbf{x})$  均以式(8.4)的连续实值函数为准。

【注 2】乍一看本式有些问题，为什么要最小化  $\ell_{\text{exp}}(\alpha_t h_t | \mathcal{D}_t)$ ？图 8.3 中的第 3 行的表达式  $h_t = \mathcal{L}(D, \mathcal{D}_t)$  不是代表着应该最小化  $\ell_{\text{exp}}(h_t | \mathcal{D}_t)$  么？或者从整体来看，第  $t$  轮迭代也应该最小化  $\ell_{\text{exp}}(H_t | \mathcal{D}) = \ell_{\text{exp}}(H_{t-1} + \alpha_t h_t | \mathcal{D})$ ，这样最终  $T$  轮迭代结束后得到的式(8.4)就可以最小化  $\ell_{\text{exp}}(H | \mathcal{D})$  了。实际上，理解了 AdaBoost 之后就会发现， $\ell_{\text{exp}}(\alpha_t h_t | \mathcal{D}_t)$  与  $\ell_{\text{exp}}(H_t | \mathcal{D})$  是等价的，详见后面的“AdaBoost 的个人推导”。另外， $h_t = \mathcal{L}(D, \mathcal{D}_t)$  也是推导的结论之一，即后面的式(8.18)，而不是无缘无故靠直觉用  $\mathcal{L}(D, \mathcal{D}_t)$  得到  $h_t$ 。

暂时不管以上疑问，权且按作者思路推导一下：

第 1 个等号与式(8.5)的区别仅在于到底针对  $\alpha_t h_t(\mathbf{x})$  还是  $H(\mathbf{x})$ ，代入即可；

第 2 个等号是考虑到  $h_t(\mathbf{x})$  和  $f(\mathbf{x})$  均只能取  $-1$  和  $+1$  两个值，其中  $\mathbb{I}(\cdot)$  为指示函数；

第 3 个等号对中括号的两项分别求  $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t}[\cdot]$ ，而  $e^{\alpha_t}$  和  $e^{-\alpha_t}$  与  $\mathbf{x}$  无关，可以作为常数项拿到  $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t}[\cdot]$  外面，而  $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t}[\mathbb{I}(f(\mathbf{x}) = h_t(\mathbf{x}))]$  表示在数据集  $D$  上、样本权值分布为  $\mathcal{D}_t$  时  $f(\mathbf{x})$  和  $h_t(\mathbf{x})$  相等次数的期望，即  $P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) = h_t(\mathbf{x}))$ ，也就是正确率，即  $(1 - \epsilon_t)$ ；同理， $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t}[\mathbb{I}(f(\mathbf{x}) \neq h_t(\mathbf{x}))]$  表示在数据集  $D$  上、样本权值分布为  $\mathcal{D}_t$  时  $f(\mathbf{x})$  和  $h_t(\mathbf{x})$  不相等次数的期望，即  $P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) \neq h_t(\mathbf{x}))$ ，也就是错误率  $\epsilon_t$ ；

第4个等号即为将 $P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) = h_t(\mathbf{x}))$ 替换为 $(1 - \epsilon_t)$ 、将 $P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) \neq h_t(\mathbf{x}))$ 替换为 $\epsilon_t$ 的结果。

## 7、式(8.12)的推导[?]

【注】本式确实有问题，虽然并不影响后面式(8.18)和式(8.19)的结果。AdaBoost 第 $t$ 轮迭代应该求解如下优化问题从而得到 $\alpha_t$ 和 $h_t(\mathbf{x})$ ：

$$(\alpha_t, h_t(\mathbf{x})) = \arg \min_{\alpha, h} \ell_{\text{exp}}(H_{t-1} + \alpha h \mid \mathcal{D})$$

对于该问题，先对于固定的任意 $\alpha > 0$ ，求解 $h_t(\mathbf{x})$ ；得到 $h_t(\mathbf{x})$ 后再求 $\alpha_t$ 。

原始文献[Friedman J H, Hastie T, Tibshirani R, et al. [Additive logistic regression : A statistical view of boosting](#)[J]. Annals of Statistics, 2000, 28(2): 337-407.]第10页(第346页)对西瓜书中式(8.12)到式(8.12)的相关推导如下：

**RESULT 1.** *The Discrete AdaBoost algorithm (population version) builds an additive logistic regression model via Newton-like updates for minimizing  $E(e^{-yF(x)})$ .*

**PROOF.** Let  $J(F) = E[e^{-yF(x)}]$ . Suppose we have a current estimate  $F(x)$  and seek an improved estimate  $F(x) + cf(x)$ . For fixed  $c$  (and  $x$ ), we expand  $J(F(x) + cf(x))$  to second order about  $f(x) = 0$ ,

$$\begin{aligned} J(F + cf) &= E[e^{-y(F(x)+cf(x))}] \\ &\approx E[e^{-yF(x)}(1 - ycf(x) + c^2 y^2 f(x)^2/2)] \\ &= E[e^{-yF(x)}(1 - ycf(x) + c^2/2)], \end{aligned}$$

since  $y^2 = 1$  and  $f(x)^2 = 1$ . Minimizing pointwise with respect to  $f(x) \in \{-1, 1\}$ , we write

$$(16) \quad f(x) = \arg \min_f E_w(1 - ycf(x) + c^2/2|x).$$

可以发现原文献中保留了参数 $c$ （即 $\alpha$ ）。当然，对于任意 $\alpha > 0$ ，并不影响推导结果。

且不管以上问题，权且按作者思路推导一下：

本式较为简单，只要将 $H_t(\mathbf{x}) = H_{t-1}(\mathbf{x}) + h_t(\mathbf{x})$ 替换式(8.5)中的 $H(\mathbf{x})$ 即可。

## 8、式(8.13)的推导

复习一下泰勒展开式：

$$f(x) = \frac{f(x_0)}{0!} + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots + \frac{f^{(n)}(x_0)}{n!}(x - x_0)^n + \dots$$

当 $f(x) = e^{-x}$ 时，在 $x_0 = 0$ 处泰勒展开：

$$e^{-x} = 1 - x + \frac{1}{2}x^2 + \dots$$

将 $x$ 换为 $f(\mathbf{x})h_t(\mathbf{x})$ ，对 $e^{-f(\mathbf{x})h_t(\mathbf{x})}$ 泰勒展开，代入即得式(8.13)。

实际上，此处保留一阶泰勒展开项即可，后面提到的 Gradient Boosting 理论框架就是只使用了一阶泰勒展开；当然二阶项为常数，也并不影响推导结果，原文献中也保留了二阶项。

## 9、式(8.14)的推导[?]

本式基于式(8.12)和式(8.13)的结果。且不管以上问题，权且按作者思路推导一下：

第 2 个等号就是将式(8.13)代入；

第 3 个等号是因为

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} \left( 1 - f(\mathbf{x})h(\mathbf{x}) + \frac{1}{2} \right) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \frac{3}{2} e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} - e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} f(\mathbf{x})h(\mathbf{x}) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \frac{3}{2} e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} \right] - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} f(\mathbf{x})h(\mathbf{x})] \end{aligned}$$

本式自变量为 $h(\mathbf{x})$ ，而 $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \frac{3}{2} e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} \right]$ 与 $h(\mathbf{x})$ 无关，也就是一个常数；只需最大化第二项 $-\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} f(\mathbf{x})h(\mathbf{x})]$ 即可，将负号去掉，原最小化问题变为最大化问题；

第 4 个等号仍是因为 $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]$ 是与自变量 $h(\mathbf{x})$ 无关的正常数（因为指数函数 $e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}$ 大于零，从而求期望 $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\cdot]$ 的结果也大于零），而最大化问题乘以某个正常数与原问题等价（例如 $\arg \max_x (1 - x^2)$ 与 $\arg \max_x 2(1 - x^2)$ 的结果均为 $x = 0$ ）；

## 10、式(8.16)的推导

类似于式(8.6)的推导中：

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H(\mathbf{x})}] = \sum_{i=1}^{|D|} \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i)H(\mathbf{x}_i)}$$

并结合式(8.15)的定义：

$$\mathcal{D}_t(\mathbf{x}_i) = \mathcal{D}(\mathbf{x}_i) \frac{e^{-f(\mathbf{x}_i)H_{t-1}(\mathbf{x}_i)}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]}$$

其中 $\mathbf{x}_i$ 表示数据集 $D$ 的第 $i$ 个样本，故与式(8.15)形式一样。针对式(8.14)最后一行的表达式：

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \frac{e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]} f(\mathbf{x})h(\mathbf{x}) \right] \\ &= \sum_{i=1}^{|D|} \mathcal{D}(\mathbf{x}_i) \frac{e^{-f(\mathbf{x}_i)H_{t-1}(\mathbf{x}_i)}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]} f(\mathbf{x}_i)h(\mathbf{x}_i) \\ &= \sum_{i=1}^{|D|} \mathcal{D}_t(\mathbf{x}_i) f(\mathbf{x}_i)h(\mathbf{x}_i) \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [f(\mathbf{x})h(\mathbf{x})] \end{aligned}$$

## 11、式(8.17)的推导

当 $f(\mathbf{x}) = h(\mathbf{x})$ 时， $f(\mathbf{x})h(\mathbf{x}) = 1$ ；

当 $f(\mathbf{x}) \neq h(\mathbf{x})$ 时， $f(\mathbf{x})h(\mathbf{x}) = -1$ ；

当 $f(\mathbf{x}) = h(\mathbf{x})$ 时， $\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x})) = 0$ ，则 $1 - 2\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x})) = 1$ ；

当 $f(\mathbf{x}) \neq h(\mathbf{x})$ 时， $\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x})) = 1$ ，则 $1 - 2\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x})) = -1$ ；

综上所述，等式两边相等。

## 12、式(8.18)的推导

本式基于式(8.17)的恒等关系，由式(8.16)推导而来。

$$\begin{aligned}
\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [f(\mathbf{x})h(\mathbf{x})] &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [1 - 2\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x}))] \\
&= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [1] - 2\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x}))] \\
&= 1 - 2\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x}))]
\end{aligned}$$

类似于式(8.14)的第 3 个和第 4 个等号，由式(8.16)的结果开始推导：

$$\begin{aligned}
h_t(\mathbf{x}) &= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [f(\mathbf{x})h(\mathbf{x})] \\
&= \arg \max_h (1 - 2\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x}))]) \\
&= \arg \max_h (-2\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x}))]) \\
&= \arg \min_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x}))]
\end{aligned}$$

此式表示理想的 $h_t(\mathbf{x})$ 在分布 $\mathcal{D}_t$ 下最小化分类误差，因此有图 8.3 第 3 行 $h_t(\mathbf{x}) = \mathcal{L}(D, \mathcal{D}_t)$ ，即分类器 $h_t(\mathbf{x})$ 可以基于分布 $\mathcal{D}_t$ 从数据集 $D$ 中训练而得，而我们在训练分类器时，一般来说最小化的损失函数就是分类误差。

### 13、式(8.19)的推导

第 1 个等号是将式(8.15)中的 $t$ 换为 $t + 1$ （同时 $t - 1$ 换为 $t$ ）；

第 2 个等号是将 $H_t(\mathbf{x}) = H_{t-1}(\mathbf{x}) + \alpha_t h_t(\mathbf{x})$ 代入分子即可；

第 3 个等号是乘以 $\frac{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_t(\mathbf{x})}]}$ 后，凑出式(8.15)的 $\mathcal{D}_t(\mathbf{x})$ 表达式，以符号 $\mathcal{D}_t(\mathbf{x})$ 替

换即得。到此之后，得到 $\mathcal{D}_{t+1}(\mathbf{x})$ 与 $\mathcal{D}_t(\mathbf{x})$ 的关系，但为了确保 $\mathcal{D}_{t+1}(\mathbf{x})$ 是一个分布，需要对得到的 $\mathcal{D}_{t+1}(\mathbf{x})$ 进行规范化，即图 8.3 第 7 行的 $Z_t$ 。式(8.19)第 3 行最后一个分式将在规范化过程被吸收。

### 14、AdaBoost 的个人推导

以上推导感觉有点乱，个人综合各种资料，个人版 AdaBoost 推导如下：

AdaBoost 的目标是学得 $T$ 个 $h_t(\mathbf{x})$ 和相应的 $T$ 个 $\alpha_t$ ，得到式(8.4)的 $H(\mathbf{x})$ ，使式(8.5)指数损失函数 $\ell_{\text{exp}}(H | \mathcal{D})$ 最小，这就是求解所谓的“**加性模型**”。特别强调一下，分类器 $h_t(\mathbf{x})$ 如何得到及其相应的权重 $\alpha_t$ 等于多少都是需要求解的（ $h_t(\mathbf{x}) = \mathcal{L}(D, \mathcal{D}_t)$ ，即基于分布 $\mathcal{D}_t$ 从数据集 $D$ 中经过最小化训练误差训练出分类器 $h_t$ ，也就是式(8.18)， $\alpha_t$ 参见式(8.18)）。

“通常这是一个复杂的优化问题（同时学得 $T$ 个 $h_t(\mathbf{x})$ 和相应的 $T$ 个 $\alpha_t$ 很困难）。**前向分步算法**求解这一优化问题的想法是：因为学习的是加法模型，如果能够从前向后，每一步只学习一个基函数 $h_t(\mathbf{x})$ 及其系数 $\alpha_t$ ，逐步逼近最小化指数损失函数 $\ell_{\text{exp}}(H | \mathcal{D})$ ，那么就可以简化优化的复杂度。”（摘自李航《统计学习方法》第 144 页，略有改动）

因此，AdaBoost 每轮迭代只需要得到一个基分类器及其投票权重，设第 $t$ 轮迭代需得到基分类器 $h_t(\mathbf{x})$ ，对应的投票权重为 $\alpha_t$ ，则集成分类器 $H_t(\mathbf{x}) = H_{t-1}(\mathbf{x}) + \alpha_t h_t(\mathbf{x})$ ，其中 $H_0(\mathbf{x}) = 0$ 。为表达式简洁，常常将 $h_t(\mathbf{x})$ 简写为 $h_t$ ， $H_t(\mathbf{x})$ 简写为 $H_t$ 。则第 $t$ 轮实际为如下优化问题（本节式(8.4)到式(8.8)已经证明了指数损失函数是分类任务原本 0/1 损失函数的一致替代损失函数）：

$$(\alpha_t, h_t) = \arg \min_{\alpha, h} \ell_{\text{exp}}(H_{t-1} + \alpha h | \mathcal{D})$$

表示每轮得到的基分类器 $h_t(\mathbf{x})$ 和对应的权重 $\alpha_t$ 是最小化集成分类器 $H_t = H_{t-1} + \alpha_t h_t$ 在数据集 $D$ 上、样本权值分布为 $\mathcal{D}$ （即初始化样本权值分布，也就是 $\mathcal{D}_1$ ）时的指数损失函数

$\ell_{\exp}(H_{t-1} + \alpha h \mid \mathcal{D})$ 的结果。这就是**前向分步算法**求解**加性模型**的思路。

根据式(8.5)将指数损失函数表达式代入，则

$$\begin{aligned}\ell_{\exp}(H_{t-1} + \alpha h \mid \mathcal{D}) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})(H_{t-1}(\mathbf{x}) + \alpha h(\mathbf{x}))}] \\ &= \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i)(H_{t-1}(\mathbf{x}_i) + \alpha h(\mathbf{x}_i))} \\ &= \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i)H_{t-1}(\mathbf{x}_i)} e^{-f(\mathbf{x}_i)\alpha h(\mathbf{x}_i)} \\ &= \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i)H_{t-1}(\mathbf{x}_i)} (e^{-\alpha} \mathbb{I}(f(\mathbf{x}_i) = h(\mathbf{x}_i)) + e^{\alpha} \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i)))\end{aligned}$$

上式推导中，由于 $f(\mathbf{x}_i)$ 和 $h(\mathbf{x}_i)$ 均只能取 $-1, +1$ 两个值，因此当 $f(\mathbf{x}_i) = h(\mathbf{x}_i)$ 时， $f(\mathbf{x}_i)h(\mathbf{x}_i) = 1$ ，当 $f(\mathbf{x}_i) \neq h(\mathbf{x}_i)$ 时， $f(\mathbf{x}_i)h(\mathbf{x}_i) = -1$ 。另外， $f(\mathbf{x}_i)$ 和 $h(\mathbf{x}_i)$ 要么相等，要么不相等，二者只能有一个为真，因此以下等式恒成立：

$$\mathbb{I}(f(\mathbf{x}_i) = h(\mathbf{x}_i)) + \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i)) = 1$$

所以

$$\begin{aligned}& e^{-\alpha} \mathbb{I}(f(\mathbf{x}_i) = h(\mathbf{x}_i)) + e^{\alpha} \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i)) \\ &= e^{-\alpha} \mathbb{I}(f(\mathbf{x}_i) = h(\mathbf{x}_i)) + e^{-\alpha} \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i)) - e^{-\alpha} \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i)) + e^{\alpha} \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i)) \\ &= e^{-\alpha} (\mathbb{I}(f(\mathbf{x}_i) = h(\mathbf{x}_i)) + \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i))) + (e^{\alpha} - e^{-\alpha}) \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i)) \\ &= e^{-\alpha} + (e^{\alpha} - e^{-\alpha}) \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i))\end{aligned}$$

将此结果代入 $\ell_{\exp}(H_{t-1} + \alpha h \mid \mathcal{D})$ ，得（注：以下表达式后面求解权重 $\alpha_t$ 时仍会使用）

$$\begin{aligned}\ell_{\exp}(H_{t-1} + \alpha h \mid \mathcal{D}) &= \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i)H_{t-1}(\mathbf{x}_i)} (e^{-\alpha} + (e^{\alpha} - e^{-\alpha}) \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i))) \\ &= \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i)H_{t-1}(\mathbf{x}_i)} e^{-\alpha} + \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i)H_{t-1}(\mathbf{x}_i)} (e^{\alpha} - e^{-\alpha}) \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i)) \\ &= e^{-\alpha} \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}'_t(\mathbf{x}_i) + (e^{\alpha} - e^{-\alpha}) \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}'_t(\mathbf{x}_i) \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i))\end{aligned}$$

上式中， $\mathcal{D}'_t(\mathbf{x}_i) = \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i)H_{t-1}(\mathbf{x}_i)}$ ， $e^{\alpha}$ 和 $e^{-\alpha}$ 与求和变量 $i$ 无关，故可以拿到求和号外面；第一项 $e^{-\alpha} \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}'_t(\mathbf{x}_i)$ 与 $h(\mathbf{x})$ 无关，因此对于任意 $\alpha > 0$ ，使 $\ell_{\exp}(H_{t-1} + \alpha h \mid \mathcal{D})$ 最小的 $h(\mathbf{x})$ 只需要使第二项最小即可，即

$$h_t = \arg \min_h (e^{\alpha} - e^{-\alpha}) \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}'_t(\mathbf{x}_i) \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i))$$

对于任意 $\alpha > 0$ ，有 $e^{\alpha} - e^{-\alpha} > 0$ ，所以上式中与 $h(\mathbf{x})$ 无关的正系数可以省略：

$$h_t = \arg \min_h \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}'_t(\mathbf{x}_i) \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i))$$

此即式(8.18)另一种表达形式。注意，为了确保 $\mathcal{D}'_t(\mathbf{x})$ 是一个分布，需要对其进行规范化，即 $\mathcal{D}_t(\mathbf{x}) = \frac{\mathcal{D}'_t(\mathbf{x})}{Z_t}$ ，然而规范化因子 $Z_t = \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}'_t(\mathbf{x}_i)$ 为常数，并不影响最小化的求解。正是基于此结论，AdaBoost通过 $h_t = \mathcal{L}(D, \mathcal{D}_t)$ 得到第 $t$ 轮的基分类器（图8.3的第3行）。

由 $\mathcal{D}_t(\mathbf{x}_i) = \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i)H_{t-1}(\mathbf{x}_i)}$ 可知（规范化只是除以常数项）

$$\begin{aligned}
 \mathcal{D}_{t+1}(\mathbf{x}_i) &= \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i) H_t(\mathbf{x}_i)} \\
 &= \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i)(H_{t-1}(\mathbf{x}_i) + \alpha_t h_t(\mathbf{x}_i))} \\
 &= \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i) H_{t-1}(\mathbf{x}_i)} e^{-f(\mathbf{x}_i) \alpha_t h_t(\mathbf{x}_i)} \\
 &= \mathcal{D}_t(\mathbf{x}_i) e^{-f(\mathbf{x}_i) \alpha_t h_t(\mathbf{x}_i)}
 \end{aligned}$$

此即类似式(8.19)的分布权重更新公式。

现在只差权重 $\alpha_t$ 表达式待求。对指数损失函数 $\ell_{\text{exp}}(H_{t-1} + \alpha h_t \mid \mathcal{D})$ 求导，得

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

令导数等于零，得

$$\begin{aligned}
 \frac{e^{-\alpha}}{e^\alpha + e^{-\alpha}} &= \frac{\sum_{i=1}^{|D|} \mathcal{D}'_t(\mathbf{x}_i) \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i))}{\sum_{i=1}^{|D|} \mathcal{D}'_t(\mathbf{x}_i)} = \sum_{i=1}^{|D|} \frac{\mathcal{D}'_t(\mathbf{x}_i)}{Z_t} \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i)) \\
 &= \sum_{i=1}^{|D|} \mathcal{D}_t(\mathbf{x}_i) \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i)) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [\mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i))] \\
 &= \epsilon_t
 \end{aligned}$$

对上述等式化简，得

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

即式(8.11)。从该式可以发现，当 $\epsilon_t = 1$ 时， $\alpha_t \rightarrow \infty$ ，此时集成分类器将由基分类器 $h_t$ 决定，而这很可能是由于过拟合产生的结果，例如不剪枝决策树，如果一直分下去，一般情况下总能得到在训练集上分类误差很小甚至为 0 的分类器，但这并没有什么意义。所以一般在 AdaBoost 中使用弱分类器，如决策树桩（即单层决策树）。

另外，由以上指数损失函数 $\ell_{\text{exp}}(H_{t-1} + \alpha h \mid \mathcal{D})$ 的推导可以发现

$$\begin{aligned}
 \ell_{\text{exp}}(H_{t-1} + \alpha h \mid \mathcal{D}) &= \sum_{i=1}^{|D|} \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i) H_{t-1}(\mathbf{x}_i)} e^{-f(\mathbf{x}_i) \alpha h(\mathbf{x}_i)} \\
 &= \sum_{i=1}^{|D|} \mathcal{D}'_t(\mathbf{x}_i) e^{-f(\mathbf{x}_i) \alpha h(\mathbf{x}_i)}
 \end{aligned}$$

这与指数损失函数 $\ell_{\text{exp}}(\alpha_t h_t \mid \mathcal{D}_t)$ 的表达式基本一致：

$$\begin{aligned}
 \ell_{\text{exp}}(\alpha_t h_t \mid \mathcal{D}_t) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [e^{-f(\mathbf{x}) \alpha_t h_t(\mathbf{x})}] \\
 &= \sum_{i=1}^{|D|} \mathcal{D}_t(\mathbf{x}_i) e^{-f(\mathbf{x}_i) \alpha_t h_t(\mathbf{x}_i)}
 \end{aligned}$$

而 $\mathcal{D}'_t(\mathbf{x})$ 的规范化过程并不影响对 $\ell_{\text{exp}}(H_{t-1} + \alpha h \mid \mathcal{D})$ 求最小化操作，因此最小化式(8.9)等价于最小化 $\ell_{\text{exp}}(H_{t-1} + \alpha h \mid \mathcal{D})$ ，这就是式(8.9)的来历，故并无问题。

到此为止，就逐一完成了图 8.3 中第 3 行的 $h_t$ 的训练（并计算训练误差）、第 6 行的权重 $\alpha_t$ 计算公式以及第 7 行的分布 $\mathcal{D}_t$ 更新公式来历的理论推导。

## 15、更进一步理解权重更新公式的含义

文献[Friedman J H, Hastie T, Tibshirani R, et al. [Additive logistic regression : A statistical view of boosting](#)[J]. Annals of Statistics, 2000, 28(2): 337-407.]第 12 页(第 348 页)有如下推论:

**COROLLARY 2.** *After each update to the weights, the weighted misclassification error of the most recent weak learner is 50%.*

**PROOF.** This follows by noting that the  $c$  that minimizes  $J(F + cf)$  satisfies

$$(21) \quad \frac{\partial J(F + cf)}{\partial c} = -E[e^{-y(F(x)+cf(x))} yf(x)] = 0.$$

The result follows since  $yf(x)$  is 1 for a correct and  $-1$  for an incorrect classification.  $\square$

即 $P_{\mathbf{x} \sim \mathcal{D}_t}(h_{t-1}(\mathbf{x}) \neq f(\mathbf{x})) = 0.5$ 。用通俗的话来说就是， $h_{t-1}$ 在数据集 $D$ 上、分布为 $\mathcal{D}_t$ 时的分类误差为 **0.5**，即相当于随机猜测（最糟糕的二分类器是分类误差为 0.5，当二分类器分类误差为 1 时相当于分类误差为 0，因为将预测结果反过来用就是了）。而 $h_t$ 由式(8.18)得到

$$h_t = \arg \min_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x}))] = \arg \min_h P_{\mathbf{x} \sim \mathcal{D}_t}(h(\mathbf{x}) \neq f(\mathbf{x}))$$

即 $h_t$ 是在数据集 $D$ 上、分布为 $\mathcal{D}_t$ 时分类误差最小的分类器，因此在数据集 $D$ 上、分布为 $\mathcal{D}_t$ 时， $h_t$ 是最好的分类器，而 $h_{t-1}$ 是最差的分类器，故二者差别最大。第 8.1 节的图 8.2 形象的说明了“集成个体应‘好而不同’”，此时可以说 $h_{t-1}$ 和 $h_t$ 非常“不同”。证明如下：

对于 $h_{t-1}$ 来说，分类误差 $\epsilon_{t-1}$ 为

$$\begin{aligned} \epsilon_{t-1} &= P_{\mathbf{x} \sim \mathcal{D}_{t-1}}(h_{t-1}(\mathbf{x}) \neq f(\mathbf{x})) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{t-1}} [\mathbb{I}(h_{t-1}(\mathbf{x}) \neq f(\mathbf{x}))] \\ &= \sum_{i=1}^{|D|} \mathcal{D}_{t-1}(\mathbf{x}_i) \mathbb{I}(h_{t-1}(\mathbf{x}_i) \neq f(\mathbf{x}_i)) \\ &= \frac{\sum_{i=1}^{|D|} \mathcal{D}_{t-1}(\mathbf{x}_i) \mathbb{I}(h_{t-1}(\mathbf{x}_i) \neq f(\mathbf{x}_i))}{\sum_{i=1}^{|D|} \mathcal{D}_{t-1}(\mathbf{x}_i) \mathbb{I}(h_{t-1}(\mathbf{x}_i) = f(\mathbf{x}_i)) + \sum_{i=1}^{|D|} \mathcal{D}_{t-1}(\mathbf{x}_i) \mathbb{I}(h_{t-1}(\mathbf{x}_i) \neq f(\mathbf{x}_i))} \end{aligned}$$

在第 $t$ 轮，根据分布更新公式(8.19)或图 3 第 7 行（规范化因子 $Z_{t-1}$ 为常量）：

$$\mathcal{D}_t = \frac{\mathcal{D}_{t-1}}{Z_{t-1}} e^{-f(\mathbf{x}) \alpha_{t-1} h_{t-1}(\mathbf{x})}$$

其中根据式(8.11)，第 $t-1$ 轮的权重

$$\alpha_{t-1} = \frac{1}{2} \ln \frac{1 - \epsilon_{t-1}}{\epsilon_{t-1}} = \ln \sqrt{\frac{1 - \epsilon_{t-1}}{\epsilon_{t-1}}}$$

代入 $\mathcal{D}_t$ 的表达式，则

$$\mathcal{D}_t = \begin{cases} \frac{\mathcal{D}_{t-1}}{Z_{t-1}} \cdot \sqrt{\frac{\epsilon_{t-1}}{1 - \epsilon_{t-1}}} & , \text{if } h_{t-1}(\mathbf{x}) = f(\mathbf{x}) \\ \frac{\mathcal{D}_{t-1}}{Z_{t-1}} \cdot \sqrt{\frac{1 - \epsilon_{t-1}}{\epsilon_{t-1}}} & , \text{if } h_{t-1}(\mathbf{x}) \neq f(\mathbf{x}) \end{cases}$$

那么 $h_{t-1}$ 在数据集 $D$ 上、分布为 $\mathcal{D}_t$ 时的分类误差 $P_{\mathbf{x} \sim \mathcal{D}_t}(h_{t-1}(\mathbf{x}) \neq f(\mathbf{x}))$ 为（注意，下式第二行的分母等于 1，因为 $\mathbb{I}(h_{t-1}(\mathbf{x}) = f(\mathbf{x})) + \mathbb{I}(h_{t-1}(\mathbf{x}) \neq f(\mathbf{x})) = 1$ ）

$$\begin{aligned}
P_{\mathbf{x} \sim \mathcal{D}_t}(h_{t-1}(\mathbf{x}) \neq f(\mathbf{x})) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [\mathbb{I}(h_{t-1}(\mathbf{x}) \neq f(\mathbf{x}))] \\
&= \frac{\sum_{i=1}^{|\mathcal{D}|} \mathcal{D}_t(\mathbf{x}_i) \mathbb{I}(h_{t-1}(\mathbf{x}_i) \neq f(\mathbf{x}_i))}{\sum_{i=1}^{|\mathcal{D}|} \mathcal{D}_t(\mathbf{x}_i) \mathbb{I}(h_{t-1}(\mathbf{x}_i) = f(\mathbf{x}_i)) + \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}_t(\mathbf{x}_i) \mathbb{I}(h_{t-1}(\mathbf{x}_i) \neq f(\mathbf{x}_i))} \\
&= \frac{\sum_{i=1}^{|\mathcal{D}|} \frac{\mathcal{D}_{t-1}(\mathbf{x}_i)}{Z_{t-1}} \cdot \sqrt{\frac{1-\epsilon_{t-1}}{\epsilon_{t-1}}} \mathbb{I}(h_{t-1}(\mathbf{x}_i) \neq f(\mathbf{x}_i))}{\sum_{i=1}^{|\mathcal{D}|} \frac{\mathcal{D}_{t-1}(\mathbf{x}_i)}{Z_{t-1}} \cdot \sqrt{\frac{\epsilon_{t-1}}{1-\epsilon_{t-1}}} \mathbb{I}(h_{t-1}(\mathbf{x}_i) = f(\mathbf{x}_i)) + \sum_{i=1}^{|\mathcal{D}|} \frac{\mathcal{D}_{t-1}(\mathbf{x}_i)}{Z_{t-1}} \cdot \sqrt{\frac{1-\epsilon_{t-1}}{\epsilon_{t-1}}} \mathbb{I}(h_{t-1}(\mathbf{x}_i) \neq f(\mathbf{x}_i))} \\
&= \frac{\sqrt{\frac{1-\epsilon_{t-1}}{\epsilon_{t-1}}} \cdot \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}_{t-1}(\mathbf{x}_i) \mathbb{I}(h_{t-1}(\mathbf{x}_i) \neq f(\mathbf{x}_i))}{\sqrt{\frac{\epsilon_{t-1}}{1-\epsilon_{t-1}}} \cdot \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}_{t-1}(\mathbf{x}_i) \mathbb{I}(h_{t-1}(\mathbf{x}_i) = f(\mathbf{x}_i)) + \sqrt{\frac{1-\epsilon_{t-1}}{\epsilon_{t-1}}} \cdot \sum_{i=1}^{|\mathcal{D}|} \mathcal{D}_{t-1}(\mathbf{x}_i) \mathbb{I}(h_{t-1}(\mathbf{x}_i) \neq f(\mathbf{x}_i))} \\
&= \frac{\sqrt{\frac{1-\epsilon_{t-1}}{\epsilon_{t-1}}} \cdot \epsilon_{t-1}}{\sqrt{\frac{\epsilon_{t-1}}{1-\epsilon_{t-1}}} \cdot (1 - \epsilon_{t-1}) + \sqrt{\frac{1-\epsilon_{t-1}}{\epsilon_{t-1}}} \cdot \epsilon_{t-1}} = \frac{1}{2}
\end{aligned}$$

## 16、能够接受带权样本的基学习算法

在林轩田老师的《机器学习技法》课程第8讲提到了两种能够接受带权样本的基学习算法，分别是 SVM 和基于随机梯度下降(SGD)的对率回归：

SVM	logistic regression
$E_{\text{in}}^u \propto C \sum_{n=1}^N u_n \widehat{\text{err}}_{\text{SVM}}$ by dual QP $\Leftrightarrow$ adjusted upper bound $0 \leq \alpha_n \leq C u_n$	$E_{\text{in}}^u \propto \sum_{n=1}^N u_n \text{err}_{\text{CE}}$ by SGD $\Leftrightarrow$ sample $(\mathbf{x}_n, y_n)$ with probability proportional to $u_n$

其实原理很简单：对于 SVM 来说，针对 P130 页的优化目标式(6.29)来说，第二项为损失项，此时每个样本的损失  $\ell_{0/1}(y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1)$  直接相加，即样本权值分布为  $\mathcal{D}(\mathbf{x}_i) = \frac{1}{m}$ ，其中  $m$  为数据集  $D$  样本个数；若样本权值更新为  $\mathcal{D}_t(\mathbf{x}_i)$ ，则此时损失求和项应该变为

$$\sum_{i=1}^m m \mathcal{D}_t(\mathbf{x}_i) \cdot \ell_{0/1}(y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1)$$

若将  $\mathcal{D}(\mathbf{x}_i) = \frac{1}{m}$  替换  $\mathcal{D}_t(\mathbf{x}_i)$ ，则就是每个样本的损失  $\ell_{0/1}(y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1)$  直接相加。

如此更改后，最后推导结果影响的是式(6.39)，将由  $C = \alpha_i + \mu_i$  变为

$$C \cdot m \mathcal{D}_t(\mathbf{x}_i) = \alpha_i + \mu_i$$

进而由  $\alpha_i, \mu_i \geq 0$  导出  $0 \leq \alpha_i \leq C \cdot m \mathcal{D}_t(\mathbf{x}_i)$ 。上述林老师课程中的  $u_n = m \mathcal{D}_t(\mathbf{x}_n)$ 。

对于基于随机梯度下降(SGD)的对率回归，每次随机选择一个样本进行梯度下降，总体上的期望损失即为式(3.27)，此时每个样本被选到的概率相同，相当于  $\mathcal{D}(\mathbf{x}_i) = \frac{1}{m}$ 。若样本权值更新为  $\mathcal{D}_t(\mathbf{x}_i)$ ，则类似于 SVM，针对式(3.27)只需要给第  $i$  项乘以  $m \mathcal{D}_t(\mathbf{x}_i)$  即可，相当于每次随机梯度下降选择样本时以概率  $\mathcal{D}_t(\mathbf{x}_i)$  选择样本  $\mathbf{x}_i$  即可。

注意，这里总的损失中出现了样本个数  $m$ 。这是因为在定义损失时未求均值，若对式(6.29)的第二项和式(3.27)乘以  $\frac{1}{m}$  则可以将  $m$  抵消掉。然而常数项在最小化式(3.27)实际上并不影响什么，对于式(6.29)来说只要选择平衡参数  $C$  时选为原来的  $m$  倍即可。



当然，正如 P177 第三段中所说，“对无法接受带权样本的基学习算法，则可通过‘重采样法’来处理，即在每一轮学习中，根据样本分布对训练集重新进行采样，再用重采样而得的样本集对基学习器进行训练”。注意，“重采样法”想要得到数据分布为  $\mathcal{D}(\mathbf{x})$  的数据集，而非以  $\mathcal{D}(\mathbf{x})$  对原数据集进行 bootstrap 采样，二者区别详见下一节“Bagging 算法的解释”。

## 8.3 Bagging 与随机森林

本节内容通俗易懂，几乎不需要什么注解。

### 1、Bagging 算法的解释

如 8.3.1 节所述，Bagging 算法采用有放回采样方式，生成  $T$  个含  $m$  个训练样本的采样集（其中  $T$  为基分类器个数， $m$  训练集样本个数），由第 2 章式(2.1)可知，生成的采样集包含约 63.2% 的训练集样本。在 Matlab 中可用如下几行代码实现 bootstrap sampling 的过程：

```
idx_bag = zeros(m,1); %m is the number of training examples
for ii=1:m
    idx_bag(ii)=randi(m);%bootstrap sampling
end
idx_oob = setdiff(1:m,idx_bag);
```

以上得到了包内样本的索引 `idx_bag`，只需按此索引从训练集中得到与训练集样本数量相同的包内样本集合即可，剩余样本索引存在 `idx_oob` 中，可以发现约有  $0.368*m$  个样本；若要训练  $T$  个基分类器，只需重复以上过程  $T$  次即可。

实际使用中，个人经常采用类似 2.2.1 节的留出法，直接随机选出训练集 67% 的样本（三分之二）用于训练，在 Matlab 中可采用如下命令完成（实际已不再是 bootstrap）：

```
n_rand = randperm(m);%m is the number of training examples
idx_train = n_rand(1:round(m*0.67));
idx_test = setdiff(1:m, idx_train);
```

以上得到了采样集的索引 `idx_train`，只需按此索引从训练集中得到采样集即可，剩余样本索引存在 `idx_test` 中；若要训练  $T$  个基分类器，只需重复以上过程  $T$  次即可。

以上实现简单直接，易于理解采样过程。实际上 Matlab 自带 `randsample` 函数可以实现以上两种过程。对于标准 Bagging 使用的 bootstrap sampling：

```
w = ones(1,m);%weights
idx_bag = randsample(1:m,m,'true',w);
idx_oob = setdiff(1:m,idx_bag);
```

第 2 行表示从整数集合  $1 \sim m$ （第 1 个参数  $1:m$ ）中以对应权重  $w$ （第 4 个参数  $w$ ），用有放回采样方式（第 3 个参数 `'true'`）随机选取  $m$  个样本（第 2 个参数  $m$ ）；此处权重  $w$  不用规范化，Matlab 在使用时会自己除以 `sum(w)` 从而将权重变为一个概率分布；`randsample` 可以实现针对每个样本以不同的权重进行采样，这一点使用起来很方便。

对于第二种以无放回采样方式选取训练集部分样本（类似留出法）：

```
idx_train = randsample(1:m,round(m*0.67));
idx_test = setdiff(1:m,idx_train);
```

第 1 行表示从整数集合  $1 \sim m$ （第 1 个参数  $1:m$ ）中以均匀分布，用无放回采样方式（无第 3 个参数，即默认为 `'false'`）随机选取 `round(m*0.67)` 个样本（第 2 个参数 `round(m*0.67)`）；`randsample` 不支持带权的无放回采样方式，即每个样本被采样的概率相等。

以上两种采样方法是有区别的，分析如下：

第 1 种标准 Bagging 使用的 bootstrap 采样使用原来约 63.2% 的样本，但采样集样本总数量保持与原样本集相同，也就是说必然有些样本在采样集中重复出现多次，因此正如西瓜书第 27 页 2.2.3 节最后一行所说，“自助法产生的数据集改变了初始数据集的分布”，即若以数据分布  $\mathcal{D}(\mathbf{x})$ （即上述使用 `randsample` 实现时的参数  $w$ ）对数据集  $D$  进行 bootstrap 采样得到包含样本个数仍为  $|D|$  的采样集  $D'$ ，则  $D'$  的数据分布一定不再是  $\mathcal{D}(\mathbf{x})$ 。在图 8.5 中，伪代码第 2 行的  $\mathcal{D}_{bs}$  可简单理解为每个样本在当前采样集中出现的概率，按 Bootstrap sampling 方法，初始训练集中约有 63.2% 的样本出现在采样集中，注意采样集与初始训练集样本个数相同，也就是说初始训练集中的有些样本在采样集中出现多次，有些样本在采样集中出现一次，有些样本在采样集中没有出现，初始训练集中每个样本在采样集中出现的频次除以样本个数  $m$  即组成  $\mathcal{D}_{bs}$ ，即自助采样产生的样本分布。

第 2 种类似的留出法只使用原数据集的部分样本，且每个样本只出现一次，即在采样集中每个样本权重相同。虽然以上无论是基于 `randperm` 函数的实现还是 `randsample` 函数的实现，都是以均匀分布从原数据集中取出约 67% 的样本，但即使采用一些特殊手段以非均匀分布从原数据集中取出约 67% 的样本，到了采样集中每个样本仍只出现一次，若使用当前采样集直接训练学习器，仍是相当于赋予采样集中每个样本相同的权重。

最后，特别提一下在实现 AdaBoost 时，“对无法接受带权样本的基学习算法”，“通过‘重采样法’来处理”时使用的“重采样法”。注意图 8.3 第 3 行的  $h_t = \mathcal{L}(D, \mathcal{D}_t)$  是“基于分布  $\mathcal{D}_t$  从数据集  $D$  中训练出分类器  $h_t$ ”，也就是说“重采样法”需要从数据集  $D$  中得到一个数据分布为  $\mathcal{D}_t$  的采样集  $D'$ （这是“重采样法”的目标），而不是以数据分布  $\mathcal{D}_t$  对数据集  $D$  进行 bootstrap 采样或其它采样。若想实现“重采样法”的目标，可简单按如下例子方法实现：

例：数据集  $D$  包含 5 个样本  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5$ ，想通过“重采样法”实现数据分布  $D = [0.1, 0.2, 0.3, 0.2, 0.2]$  的采样集  $D'$ ，则

$$D' = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_3, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_5\}$$

即以数据分布为比例简单扩充数据集  $D$  得到采样集  $D'$ 。

另外，以上分析和代码实现只是个人看法和偏好，仅供参考，勿被误导^\_^

## 2、式(8.20)的解释

仅注意一点即可，每个基分类器对应的 36.8% 的剩余样本均是不同的（可能会有样本出现在多个基分类器对应的 36.8% 的剩余样本集合中），因此总体考虑包外估计时，整体考虑的样本肯定大于训练集 36.8% 的样本个数，但这些样本并不一定是所有  $T$  个基分类器的包外样本，因此在式(8.20)投票时，针对每个样本  $\mathbf{x}$  不会有  $T$  票出现（ $\mathbb{I}(\mathbf{x} \notin D_t)$  一般不会  $T$  次均为真），即“仅考虑那些未使用  $\mathbf{x}$  训练的基学习器在  $\mathbf{x}$  上的预测”。

## 3、式(8.21)的推导

本式直接使用式(8.20)的包外预测结果。但本式直接除以  $|D|$ （训练集  $D$  样本个数），也就是说此处假设  $T$  个基分类器的各自的包外样本的并集一定为训练集  $D$ 。实际上，这个事实成立的概率也是比较大的，可以计算一下：样本属于包内的概率为 0.632，那么  $T$  次独立的随机采样均属于包内的概率为  $0.632^T$ ，当  $T = 5$  时， $0.632^T \approx 0.1$ ，当  $T = 10$  时， $0.632^T \approx 0.01$ ，这么来看的话  $T$  个基分类器的各自的包外样本的并集为训练集  $D$  的概率的确挺大的。

## 4、随机森林的解释

在 8.3.2 节开篇第一句话就解释了随机森林的概念：随机森林是 Bagging 的一个扩展变体，是以决策树为基学习器构建 Bagging 集成的基础上，进一步在决策树的训练过程中引入

了随机属性选择。

完整版随机森林当然更复杂，这时只须知道两个重点：(1) 以决策树为基学习器；(2) 在基学习器训练过程中，选择划分属性时只使用当前结点属性集合的一个子集。

## 8.4 结合策略

本节内容通俗易懂，几乎不需要什么注解。

### 1、硬投票和软投票的解释

书中第 183 页提到了硬投票(hard voting)和软投票(soft voting)，本页左侧注释也提到多数投票法的英文术语使用不太一致，有文献称为 majority voting。本人看到有些文献中，硬投票使用 majority voting（多数投票），软投票使用 probability voting（概率投票），所以还是具体问题具体分析比较稳妥。

### 2、元学习器(meta-learner)的解释

书中第 183 页最后一行提到了元学习器(meta-learner)，简单解释一下，因为理解 meta 的含义有时对于理解论文中的核心思想很有帮助。

元(meta)，非常抽象，例如此处的含义，即次级学习器，或者说基于学习器结果的学习器；另外还有元语言，就是描述计算机语言的语言，还有元数学，研究数学的数学等等；

另外，论文中经常出现的还有 meta-strategy，即元策略或元方法，比如说你的研究问题是多分类问题，那么你提出了一种方法，例如对输入特征进行变换（或对输出类别做某种变换），然后再基于普通的多分类方法进行预测，这时你的方法可以看成是一种通用的框架，它虽然针对多分类问题开发，但它需要某个具体多分类方法配合才能实现，那么这样的方法是一种更高层级的方法，可以称为是一种 meta-strategy。

不太理解没关系，等看论文中碰到了再回来看看 stacking 算法，细细品味，定能悟道。

### 3、Stacking 算法的解释

该算法其实非常简单，对于数据集，试想你现在有了  $T$  个基分类器预测结果，也就是说数据集中的每个样本均有  $T$  个预测结果，那么怎么结合这  $T$  个预测结果呢？

本节名为“结合策略”，告诉你各种结合方法，但其实最简单的方法就是基于这  $T$  个预测结果再进行一次学习，即针对每个样本，将这  $T$  个预测结果作为输入特征，类别仍为原来的类别，既然无法抉择如何将这  $T$  个结果进行结合，那么就“学习”一下吧。

图 8.9 伪代码第 9 行中  $((z_{i1}, z_{i2}, \dots, z_{iT}), y_i)$  将第  $i$  个样本进行变换，特征为  $T$  个基学习器的输出  $(z_{i1}, z_{i2}, \dots, z_{iT})$ ，类别标记仍为原来的  $y_i$ ，将所有训练集  $D$  中的样本进行转换得到新的数据集  $D'$  后，再基于  $D'$  进行一次学习即可，也就是 Stacking 算法。

至于说图 8.9 中伪代码第 1 行到第 3 行使用的数据集与第 5 行到第 10 行使用的数据集  $D$  之间的关系，在图 8.9 下方的一段话有详细的讨论，不再赘述。

## 8.5 多样性

若无特别需求，本节内容只要了解 8.5.3 节介绍的有关多样性增强的一些技巧即可。

### 1、式(8.28)的解释

该式本身并无太多疑问，但该式上方一行话“则集成的‘分歧’是”是令人困惑的，为

什么集成的“分歧”是式(8.28)呢？集成不就是 $H(\mathbf{x})$ 么？

个人感觉此处更应该写为“则个体学习器‘分歧’的加权均值是”，此表述在式(8.36)上方解释 $\bar{A}$ 时被使用，否则可以会有点令人费解。

## 2、式(8.31)与式(8.28)的等价性推导

式(8.31)与式(8.28)均为 $\bar{A}(h | \mathbf{x})$ ，因此应该是等价的，接下来进行推导。

(1)首先，化简式(8.31)：

$$\begin{aligned}
 \bar{A}(h | \mathbf{x}) &= \sum_{i=1}^T w_i E(h_i | \mathbf{x}) - E(H | \mathbf{x}) \\
 &= \sum_{i=1}^T w_i [f(\mathbf{x}) - h_i(\mathbf{x})]^2 - [f(\mathbf{x}) - H(\mathbf{x})]^2 \\
 &= \sum_{i=1}^T [w_i (f(\mathbf{x}))^2 - 2w_i f(\mathbf{x}) h_i(\mathbf{x}) + w_i (h_i(\mathbf{x}))^2] \\
 &\quad - [(f(\mathbf{x}))^2 - 2f(\mathbf{x})H(\mathbf{x}) + (H(\mathbf{x}))^2] \\
 &= (f(\mathbf{x}))^2 \sum_{i=1}^T w_i - 2f(\mathbf{x}) \sum_{i=1}^T w_i h_i(\mathbf{x}) + \sum_{i=1}^T w_i (h_i(\mathbf{x}))^2 \\
 &\quad - (f(\mathbf{x}))^2 + 2f(\mathbf{x})H(\mathbf{x}) - (H(\mathbf{x}))^2
 \end{aligned}$$

根据式(8.23)定义 $H(\mathbf{x}) = \sum_{i=1}^T w_i h_i(\mathbf{x})$ ，注意约束条件 $\sum_{i=1}^T w_i = 1$ ，因此

$$(f(\mathbf{x}))^2 \sum_{i=1}^T w_i = (f(\mathbf{x}))^2, \quad f(\mathbf{x}) \sum_{i=1}^T w_i h_i(\mathbf{x}) = f(\mathbf{x})H(\mathbf{x})$$

代入上式化简结果，得

$$\bar{A}(h | \mathbf{x}) = \sum_{i=1}^T w_i (h_i(\mathbf{x}))^2 - (H(\mathbf{x}))^2$$

(2)接下来，化简式(8.28)：

$$\begin{aligned}
 \bar{A}(h | \mathbf{x}) &= \sum_{i=1}^T w_i [h_i(\mathbf{x}) - H(\mathbf{x})]^2 \\
 &= \sum_{i=1}^T w_i (h_i(\mathbf{x}))^2 - \sum_{i=1}^T 2w_i h_i(\mathbf{x}) H(\mathbf{x}) + \sum_{i=1}^T w_i (H(\mathbf{x}))^2 \\
 &= \sum_{i=1}^T w_i (h_i(\mathbf{x}))^2 - 2H(\mathbf{x}) \sum_{i=1}^T w_i h_i(\mathbf{x}) + (H(\mathbf{x}))^2 \sum_{i=1}^T w_i
 \end{aligned}$$

将 $H(\mathbf{x}) = \sum_{i=1}^T w_i h_i(\mathbf{x})$ 和 $\sum_{i=1}^T w_i = 1$ 代入上式化简结果，得

$$\begin{aligned}
 \bar{A}(h | \mathbf{x}) &= \sum_{i=1}^T w_i (h_i(\mathbf{x}))^2 - 2H(\mathbf{x})H(\mathbf{x}) + (H(\mathbf{x}))^2 \\
 &= \sum_{i=1}^T w_i (h_i(\mathbf{x}))^2 - (H(\mathbf{x}))^2
 \end{aligned}$$

可以发现，式(8.28)的化简结果与式(8.31)的化简结果相等，即两种定义等价。

### 3、式(8.32)的解释

式(8.31)的意义在于，对于示例 $\mathbf{x}$ 有 $\overline{A}(h | \mathbf{x}) = \overline{E}(h | \mathbf{x}) - E(H | \mathbf{x})$ 成立，即个体学习器分歧的加权均值等于个体学习器误差的加权均值减去集成 $H(\mathbf{x})$ 的误差。

将这个结论应用于全样本上，即为式(8.32)。

例如 $A_i = \int A(h_i | \mathbf{x})p(\mathbf{x})d\mathbf{x}$ ，这是将 $\mathbf{x}$ 作为连续变量来处理的，所以这里是概率密度 $p(\mathbf{x})$ 和积分号；若按离散变量来处理，则变为 $A_i = \sum_{\mathbf{x} \in D} A(h_i | \mathbf{x})p_{\mathbf{x}}$ ；其实高等数学中讲过，积分就是连续求和。

### 4、列联表(contingency table)的解释

第 187 页第 2 行出现了“列联表”，看论文时偶尔能遇到该术语。列联表是统计学中双变量相关分析时常见的概念。例如，现在要调查某医院各职业（医生、护士、行政人员）和性别（男、女）的相关关系，则需要先列出列联表：

		职业		
		医生	护士	行政人员
性别	男	$a$	$b$	$c$
	女	$d$	$e$	$f$

其中男医生人数为 $a$ ，男护士人数为 $b$ ，男行政人员人数为 $c$ ，女医生人数为 $d$ ，女护士人数为 $e$ ，女行政人员人数为 $f$ ，该医院总人数（或参与调查的人数）为 $a + b + c + d + e + f$ 。有了列联表，就可以计算一些统计量，如克莱姆相关系数。

列联表又称交叉资料表、交互分类表等，可参见百度百科词条[列联表](#)。

### 5、式(8.40)的解释

当 $p_1 = p_2$ 时， $\kappa = 0$ ；当 $p_1 = 1$ 时， $\kappa = 1$ ；一般来说 $p_1 \geq p_2$ ，即 $\kappa \geq 0$ ，但偶尔也有 $p_1 < p_2$ 的情况，此时 $\kappa < 0$ 。

有关 $p_1, p_2$ 的意义参见式(8.41)和式(8.42)的解释。

### 6、式(8.41)的解释

分子 $a + d$ 为分类器 $h_i$ 与 $h_j$ 在数据集 $D$ 上预测结果相同的样本数目，分母为数据集 $D$ 总样本数目，因此 $p_1$ 为两个分类器 $h_i$ 与 $h_j$ 预测结果相同的概率。

若 $a + d = m$ ，即分类器 $h_i$ 与 $h_j$ 对数据集 $D$ 所有样本预测结果均相同，此时 $p_1 = 1$ 。

### 7、式(8.42)的解释

将式(8.42)拆分为如下形式，将会很容易理解其含义：

$$p_2 = \frac{a+b}{m} \cdot \frac{a+c}{m} + \frac{c+d}{m} \cdot \frac{b+d}{m}$$

其中 $\frac{a+b}{m}$ 为分类器 $h_i$ 将样本预测为+1的概率， $\frac{a+c}{m}$ 为分类器 $h_j$ 将样本预测为+1的概率，二者相乘 $\frac{a+b}{m} \cdot \frac{a+c}{m}$ 可理解为分类器 $h_i$ 与 $h_j$ 将样本预测为+1的概率； $\frac{c+d}{m}$ 为分类器 $h_i$ 将样本预测为-1的概率， $\frac{b+d}{m}$ 为分类器 $h_j$ 将样本预测为-1的概率，二者相乘 $\frac{c+d}{m} \cdot \frac{b+d}{m}$ 可理解为分类器 $h_i$ 与 $h_j$ 将样本预测为-1的概率。

注意 $\frac{a+b}{m} \cdot \frac{a+c}{m}$ 与 $\frac{a}{m}$ 的不同， $\frac{c+d}{m} \cdot \frac{b+d}{m}$ 与 $\frac{d}{m}$ 的不同：

$$\frac{a+b}{m} \cdot \frac{a+c}{m} = p(h_i = +1)p(h_j = +1), \frac{a}{m} = p(h_i = +1, h_j = +1)$$

$$\frac{c+d}{m} \cdot \frac{b+d}{m} = p(h_i = -1)p(h_j = -1), \frac{d}{m} = p(h_i = -1, h_j = -1)$$

即  $\frac{a+b}{m}$ 、 $\frac{a+c}{m}$  和  $\frac{c+d}{m}$ 、 $\frac{b+d}{m}$  是分别考虑分类器  $h_i$  与  $h_j$  时的概率 ( $h_i$  与  $h_j$  独立), 而  $\frac{a}{m}$  和  $\frac{d}{m}$  是同时考虑  $h_i$  与  $h_j$  时的概率 (联合概率)。

## 8、多样性增强的解释

在 8.5.3 节介绍了四种多样性增强的方法, 通俗易懂, 几乎不需要什么注解, 仅强调几个概念:

(1) 数据样本扰动中提到了“不稳定基学习器”(例如决策树、神经网络等) 和“稳定基学习器”(例如线性学习器、支持向量机、朴素贝叶斯、 $k$ 近邻学习器等), 对稳定基学习器进行集成时数据样本扰动技巧效果有限。这也可以解释为什么随机森林和 GBDT 等以决策树为基分学习器的集成方法很成功吧, Gradient Boosting 和 Bagging 都是以数据样本扰动来增强多样性的; 而且, 掌握这个经验后在实际工程应用中就可以排除一些候选基分类器, 但论文中的确经常见到以支持向量机为基分类器 Bagging 实现, 这可能是由于 LIBSVM 简单易用的原因吧。

(2) 图 8.11 随机子空间算法, 针对每个基分类器  $h_t$  在训练时使用了原数据集的部分输入属性 (未必是初始属性, 详见第 189 页左上注释), 因此在最终集成时 (图 8.11 最后一行) 也要使用相同的部分属性。

(3) 输出表示扰动中提到了“翻转法”(Flipping Output), 看起来是一个很无厘头的技巧, 为什么要将训练样本的标记改变呢? 若认为原训练样本标记是完全可靠的, 这不是人为地加入噪声么? 但西瓜书作者 2017 年提出的深度森林 (参见 IJCAI'17: Deep Forest: Towards an Alternative to Deep Neural Networks) 模型中也用到了该技巧, 正如本小节名为“多样性增强”, 虽然从局部来看引入了标记噪声, 但从模型集成的角度来说却是有益的。

## 8.6、GB(Gradient Boosting)/GBDT/XGBoost

在集成学习中, 梯度提升(Gradient Boosting, GB)、梯度提升树(GB Decision Tree, GBDT) 很常见, 尤其是近几年非常流行的 XGBoost 很是耀眼, 此处单独介绍对比这些概念。

### 1、梯度下降法

本部分内容参考了[孙文瑜, 徐成贤, 朱德通. 最优化方法[M]. 2 版. 北京: 高等教育出版社, 2010: 139-140.]。

设目标函数  $f(\mathbf{x})$  在  $\mathbf{x}_k$  附近连续可微, 且  $\nabla f(\mathbf{x}_k) = \left. \frac{\nabla f(\mathbf{x})}{\nabla \mathbf{x}} \right|_{\mathbf{x}=\mathbf{x}_k} \neq 0$ 。将  $f(\mathbf{x})$  在  $\mathbf{x}_k$  处进行一阶 Taylor 展开

$$f(\mathbf{x}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top (\mathbf{x} - \mathbf{x}_k)$$

记  $\mathbf{x} - \mathbf{x}_k = \Delta \mathbf{x}$ , 则上式可写为

$$f(\mathbf{x}_k + \Delta \mathbf{x}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \Delta \mathbf{x}$$

显然, 若  $\nabla f(\mathbf{x}_k)^\top \Delta \mathbf{x} < 0$  则有  $f(\mathbf{x}_k + \Delta \mathbf{x}) < f(\mathbf{x}_k)$ , 即相比于  $f(\mathbf{x}_k)$ , 自变量增量  $\Delta \mathbf{x}$  会使  $f(\mathbf{x})$  函数值下降; 若要使  $f(\mathbf{x}) = f(\mathbf{x}_k + \Delta \mathbf{x})$  下降最快, 只要选择  $\Delta \mathbf{x}$  使  $\nabla f(\mathbf{x}_k)^\top \Delta \mathbf{x}$  最小即可, 而此时  $\nabla f(\mathbf{x}_k)^\top \Delta \mathbf{x} < 0$ , 因此使绝对值  $|\nabla f(\mathbf{x}_k)^\top \Delta \mathbf{x}|$  最大即可。将  $\Delta \mathbf{x}$  分成两

部分： $\Delta \mathbf{x} = \alpha_k \mathbf{d}_k$ ，其中 $\mathbf{d}_k$ 为待求单位向量， $\alpha_k > 0$ 为待解常量； $\mathbf{d}_k$ 表示往哪个方向改变 $\mathbf{x}$ 函数值下降最快，而 $\alpha_k$ 表示沿这个方向的步长。因此，求解 $\Delta \mathbf{x}$ 的问题变为

$$(\alpha_k, \mathbf{d}_k) = \arg \min_{\alpha, \mathbf{d}} \nabla f(\mathbf{x}_k)^\top \alpha \mathbf{d}$$

将以上优化问题分为两步求解，即

$$\mathbf{d}_k = \arg \min_{\mathbf{d}} \nabla f(\mathbf{x}_k)^\top \mathbf{d} \quad \text{s.t.} \quad \|\mathbf{d}\|_2 = 1$$

$$\alpha_k = \arg \min_{\alpha} \nabla f(\mathbf{x}_k)^\top \mathbf{d}_k \alpha$$

以上求解 $\alpha_k$ 的优化问题明显有问题，因为对于 $\nabla f(\mathbf{x}_k)^\top \mathbf{d}_k < 0$ 来说，显然 $\alpha_k = +\infty$ 时取的最小值，求解 $\alpha_k$ 应该求解如下优化问题：

$$\alpha_k = \arg \min_{\alpha} f(\mathbf{x}_k + \alpha \mathbf{d}_k)$$

对于凸函数来说，以上两步可以得到最优解；但对于非凸函数来说，联合求解得到 $\mathbf{d}_k$ 和 $\alpha_k$ ，与先求 $\mathbf{d}_k$ 然后基于此再求 $\alpha_k$ 的结果应该有时是不同的。

由 Cauchy-Schwartz 不等式

$$|\nabla f(\mathbf{x}_k)^\top \mathbf{d}_k| \leq \|\nabla f(\mathbf{x}_k)\|_2 \|\mathbf{d}_k\|_2$$

可知，当且仅当 $\mathbf{d}_k = -\frac{\nabla f(\mathbf{x}_k)}{\|\nabla f(\mathbf{x}_k)\|_2}$ 时， $\nabla f(\mathbf{x}_k)^\top \mathbf{d}_k$ 最小， $-\nabla f(\mathbf{x}_k)^\top \mathbf{d}_k$ 最大。

对于 $\alpha_k$ ，若 $f(\mathbf{x}_k + \alpha \mathbf{d}_k)$ 对 $\alpha$ 的导数存在，则可简单求解如下单变量方程即可：

$$\frac{\partial f(\mathbf{x}_k + \alpha \mathbf{d}_k)}{\partial \alpha} = 0$$

例 1：试求 $f(x) = x^2$ 在 $x_k = 2$ 处的梯度方向 $\mathbf{d}_k$ 和步长 $\alpha_k$ 。

解：对 $f(x)$ 在 $x_k = 2$ 处进行一阶 Taylor 展开：

$$\begin{aligned} f(x) &= f(x_k) + f'(x_k)(x - x_k) \\ &= x_k^2 + 2x_k(x - x_k) \\ &= x_k^2 + 2x_k \alpha d \end{aligned}$$

由于此时自变量为一维，因此只有两个方向可选，要么正方向，要么负方向。此时 $f'(x_k) = 4$ ，

因此 $\mathbf{d}_k = -\frac{f'(x_k)}{|f'(x_k)|} = -1$ 。接下来求 $\alpha_k$ ，将 $x_k$ 和 $\mathbf{d}_k$ 代入：

$$f(x_k + \alpha \mathbf{d}_k) = f(2 - \alpha) = (2 - \alpha)^2$$

进而有

$$\frac{\partial f(x_k + \alpha \mathbf{d}_k)}{\partial \alpha} = -2(2 - \alpha)$$

令导数等于 0，得 $\alpha_k = 2$ 。此时

$$\Delta x = \alpha_k \mathbf{d}_k = -2$$

则 $x_k + \Delta x = 0$ ，函数值 $f(x_k + \Delta x) = 0$ 。

例 2：试求 $f(\mathbf{x}) = \|\mathbf{x}\|_2^2 = \mathbf{x}^\top \mathbf{x}$ 在 $\mathbf{x}_k = [x_k^1, x_k^2]^\top = [3, 4]^\top$ 处的梯度方向 $\mathbf{d}_k$ 和步长 $\alpha_k$ 。

解：对 $f(\mathbf{x})$ 在 $\mathbf{x}_k = [x_k^1, x_k^2]^\top = [3, 4]^\top$ 处进行一阶 Taylor 展开：

$$\begin{aligned}
f(\mathbf{x}) &= f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top (\mathbf{x} - \mathbf{x}_k) \\
&= \|\mathbf{x}\|_2^2 + 2\mathbf{x}_k^\top (\mathbf{x} - \mathbf{x}_k) \\
&= \|\mathbf{x}\|_2^2 + 2\mathbf{x}_k^\top \alpha \mathbf{d}
\end{aligned}$$

此时  $\nabla f(\mathbf{x}_k) = [6, 8]^\top$ , 因此  $\mathbf{d}_k = -\frac{\nabla f(\mathbf{x}_k)}{\|\nabla f(\mathbf{x}_k)\|_2} = [-0.6, -0.8]^\top$ 。接下来求  $\alpha_k$ , 将  $\mathbf{x}_k$  和  $\mathbf{d}_k$  代入:

$$\begin{aligned}
f(\mathbf{x}_k + \alpha \mathbf{d}_k) &= (3 - 0.6\alpha)^2 + (4 - 0.8\alpha)^2 \\
&= \alpha^2 - 10\alpha + 25 \\
&= (\alpha - 5)^2
\end{aligned}$$

因此可得  $\alpha_k = 5$  (或对  $\alpha$  求导, 再令导数等于 0)。此时

$$\Delta \mathbf{x} = \alpha_k \mathbf{d}_k = [-3, -4]^\top$$

则  $\mathbf{x}_k + \Delta \mathbf{x} = [0, 0]^\top$ , 函数值  $f(\mathbf{x}_k + \Delta \mathbf{x}) = 0$ 。

通过以上分析, 只想强调两点:

(1) 梯度下降法求解下降最快的方向  $\mathbf{d}_k$  时应该求解如下优化问题:

$$\mathbf{d}_k = \arg \min_{\mathbf{d}} \nabla f(\mathbf{x}_k)^\top \mathbf{d} \quad \text{s.t.} \quad \|\mathbf{d}\|_2 = C$$

其中  $C$  为常量, 即不必严格限定  $\|\mathbf{d}_k\|_2 = 1$ , 只要固定向量长度, 与  $\alpha_k$  搭配即可。

(2) 梯度下降法求解步长  $\alpha_k$  应该求解如下优化问题:

$$\alpha_k = \arg \min_{\alpha} f(\mathbf{x}_k + \alpha \mathbf{d}_k)$$

实际应用中, 很多时候不会去求最优的  $\alpha_k$ , 而是靠经验设置一个步长。

## 2、从梯度下降的角度解释 AdaBoost

AdaBoost 第  $t$  轮迭代时最小化式(8.5)的指数损失函数

$$\ell_{\text{exp}}(H_t \mid \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_t(\mathbf{x})}] = \sum_{\mathbf{x} \in D} \mathcal{D}(\mathbf{x}) e^{-f(\mathbf{x})H_t(\mathbf{x})}$$

对  $\ell_{\text{exp}}(H_t \mid \mathcal{D})$  每一项在  $H_{t-1}$  处泰勒展开

$$\begin{aligned}
\ell_{\text{exp}}(H_t \mid \mathcal{D}) &\approx \sum_{\mathbf{x} \in D} \mathcal{D}(\mathbf{x}) (e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} - f(\mathbf{x})e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}(H_t(\mathbf{x}) - H_{t-1}(\mathbf{x}))) \\
&= \sum_{\mathbf{x} \in D} \mathcal{D}(\mathbf{x}) (e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} - e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}f(\mathbf{x})\alpha_t h_t(\mathbf{x})) \\
&= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} - e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}f(\mathbf{x})\alpha_t h_t(\mathbf{x})]
\end{aligned}$$

其中  $H_t = H_{t-1} + \alpha_t h_t$ 。注意:  $\alpha_t, h_t$  是第  $t$  轮待解的变量。

另外补充一下, 在上式展开中的变量为  $H_t(\mathbf{x})$ , 在  $H_{t-1}$  处一阶导数为

$$\left. \frac{\partial e^{-f(\mathbf{x})H_t(\mathbf{x})}}{\partial H_t(\mathbf{x})} \right|_{H_t(\mathbf{x})=H_{t-1}(\mathbf{x})} = -f(\mathbf{x})e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}$$

如果看不习惯上述泰勒展开过程, 可令变量  $z = H_t(\mathbf{x})$  和函数  $g(z) = e^{-f(\mathbf{x})z}$ , 对  $g(z)$  在  $z_0 = H_{t-1}(\mathbf{x})$  处泰勒展开, 得



$$\begin{aligned}
 g(z) &\approx g(z_0) + g'(z_0)(z - z_0) \\
 &= g(z_0) - f(\mathbf{x})e^{-f(\mathbf{x})z_0}(z - z_0) \\
 &= e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} - e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}f(\mathbf{x})(H_t(\mathbf{x}) - H_{t-1}(\mathbf{x})) \\
 &= e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} - e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}f(\mathbf{x})\alpha_t h_t(\mathbf{x})
 \end{aligned}$$

注意此处 $h_t(\mathbf{x}) \in \{-1, +1\}$ ，类似于梯度下降法中的约束 $\|\mathbf{d}_k\|_2 = 1$ 。

类似于梯度下降法求解下降最快的方向 $\mathbf{d}_k$ ，此处先求 $h_t$ （先不管 $\alpha_t$ ）：

$$h_t = \arg \min_h \sum_{\mathbf{x} \in D} \mathcal{D}(\mathbf{x}) (-e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}f(\mathbf{x})h(\mathbf{x})) \quad \text{s.t. } h(\mathbf{x}) \in \{-1, +1\}$$

将负号去掉，最小化变为最大化问题

$$\begin{aligned}
 h_t &= \arg \max_h \sum_{\mathbf{x} \in D} \mathcal{D}(\mathbf{x}) (e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}f(\mathbf{x})h(\mathbf{x})) \\
 &= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}f(\mathbf{x})h(\mathbf{x})] \quad \text{s.t. } h(\mathbf{x}) \in \{-1, +1\}
 \end{aligned}$$

这就是式(8.14)的第3个等号的结果，因此其余推导参见8.2节即可。

由于这里的 $h(\mathbf{x})$ 约束较强，因此不能直接取负梯度方向，书中经过推导得到了 $h_t(\mathbf{x})$ 的表达式，即式(8.18)。实际上，可以将此结果理解为满足约束条件的最快下降方向。

求得 $h_t(\mathbf{x})$ 之后再求 $\alpha_t$ （8.2节“AdaBoost的个人推导”注解中已经写过一遍，此处仅粘贴至此，具体参见8.2节注解，尤其是 $\ell_{\text{exp}}(H_{t-1} + \alpha h_t \mid \mathcal{D})$ 表达式的由来）：

$$\alpha_k = \arg \min_{\alpha} \ell_{\text{exp}}(H_{t-1} + \alpha h_t \mid \mathcal{D})$$

对指数损失函数 $\ell_{\text{exp}}(H_{t-1} + \alpha h_t \mid \mathcal{D})$ 求导，得

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

令导数等于零，得

$$\begin{aligned}
 \frac{e^{-\alpha}}{e^{\alpha} + e^{-\alpha}} &= \frac{\sum_{i=1}^{|D|} \mathcal{D}'_t(\mathbf{x}_i) \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i))}{\sum_{i=1}^{|D|} \mathcal{D}'_t(\mathbf{x}_i)} = \sum_{i=1}^{|D|} \frac{\mathcal{D}'_t(\mathbf{x}_i)}{Z_t} \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i)) \\
 &= \sum_{i=1}^{|D|} \mathcal{D}_t(\mathbf{x}_i) \mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i)) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [\mathbb{I}(f(\mathbf{x}_i) \neq h(\mathbf{x}_i))] \\
 &= \epsilon_t
 \end{aligned}$$

对上述等式化简，得

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

即式(8.11)。

通过以上推导可以发现：AdaBoost 每一轮的迭代就是基于梯度下降法求解损失函数为指数损失函数的二分类问题（约束条件 $h_t(\mathbf{x}) \in \{-1, +1\}$ ）。

### 3、梯度提升(Gradient Boosting)

将 AdaBoost 的问题一般化，即不限定损失函数为指数损失函数，也不局限于二分类问题，则可以将式(8.5)写为更一般化的形式

$$\begin{aligned}\ell(H_t | \mathcal{D}) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\text{err}(H_t(\mathbf{x}), f(\mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\text{err}(H_{t-1}(\mathbf{x}) + \alpha_t h_t(\mathbf{x}), f(\mathbf{x}))]\end{aligned}$$

易知，当  $f(\mathbf{x}) \in \{-1, +1\}$  且  $\text{err}(H_t(\mathbf{x}), f(\mathbf{x})) = e^{-f(\mathbf{x})H_t(\mathbf{x})}$  时，就是式(8.5)。当是回归问题时， $f(\mathbf{x}) \in \mathbb{R}$ ，损失函数可使用平方损失  $\text{err}(H_t(\mathbf{x}), f(\mathbf{x})) = (H_t(\mathbf{x}) - f(\mathbf{x}))^2$ 。

针对该一般化的损失函数和一般的学习问题，要通过  $T$  轮迭代得到学习器

$$H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x})$$

类似于 AdaBoost，第  $t$  轮得到  $\alpha_t, h_t(\mathbf{x})$ ，可先对损失函数在  $H_{t-1}(\mathbf{x})$  处进行泰勒展开：

$$\begin{aligned}\ell(H_t | \mathcal{D}) &\approx \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \text{err}(H_{t-1}(\mathbf{x}), f(\mathbf{x})) + \left. \frac{\partial \text{err}(H_t(\mathbf{x}), f(\mathbf{x}))}{\partial H_t(\mathbf{x})} \right|_{H_t(\mathbf{x})=H_{t-1}(\mathbf{x})} (H_t(\mathbf{x}) - H_{t-1}(\mathbf{x})) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \text{err}(H_{t-1}(\mathbf{x}), f(\mathbf{x})) + \left. \frac{\partial \text{err}(H_t(\mathbf{x}), f(\mathbf{x}))}{\partial H_t(\mathbf{x})} \right|_{H_t(\mathbf{x})=H_{t-1}(\mathbf{x})} \alpha_t h_t(\mathbf{x}) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\text{err}(H_{t-1}(\mathbf{x}), f(\mathbf{x}))] + \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \left. \frac{\partial \text{err}(H_t(\mathbf{x}), f(\mathbf{x}))}{\partial H_t(\mathbf{x})} \right|_{H_t(\mathbf{x})=H_{t-1}(\mathbf{x})} \alpha_t h_t(\mathbf{x}) \right]\end{aligned}$$

注意，在上式展开中的变量为  $H_t(\mathbf{x})$ ，且有  $H_t(\mathbf{x}) = H_{t-1}(\mathbf{x}) + \alpha_t h_t(\mathbf{x})$ （类似于梯度下降法中  $\mathbf{x} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$ ）。上式中括号内第 1 项为常量  $\ell(H_{t-1} | \mathcal{D})$ ，最小化  $\ell(H_t | \mathcal{D})$  只须最小化第 2 项即可。先不考虑权重  $\alpha_t$ ，求解如下优化问题可得  $h_t(\mathbf{x})$ ：

$$h_t(\mathbf{x}) = \arg \min_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \left. \frac{\partial \text{err}(H_t(\mathbf{x}), f(\mathbf{x}))}{\partial H_t(\mathbf{x})} \right|_{H_t(\mathbf{x})=H_{t-1}(\mathbf{x})} h(\mathbf{x}) \right] \quad \text{s.t. constraints for } h(\mathbf{x})$$

解得  $h_t(\mathbf{x})$  之后，再求解如下优化问题可得权重  $\alpha_t$ ：

$$\alpha_t = \arg \min_{\alpha} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\text{err}(H_{t-1}(\mathbf{x}) + \alpha h_t(\mathbf{x}), f(\mathbf{x}))]$$

以上就是梯度提升(Gradient Boosting)的理论框架，即每轮通过梯度(Gradient)下降的方式将  $T$  个弱学习器提升(Boosting)为强学习器。可以看出 AdaBoost 是其特殊形式。

Gradient Boosting 算法的官方版本参见[Friedman J H. [Greedy function approximation: a gradient boosting machine](#)[J]. Annals of statistics, 2001: 1189-1232.]的第 5-6 页(第 1193-1194 页)：

#### ALGORITHM 1 (Gradient Boost).

1.  $F_0(\mathbf{x}) = \arg \min_{\rho} \sum_{i=1}^N L(y_i, \rho)$
2. For  $m = 1$  to  $M$  do:
3.  $\tilde{y}_i = -\left[ \frac{\partial L(y_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)} \right]_{F(\mathbf{x})=F_{m-1}(\mathbf{x})}, i = 1, N$
4.  $\mathbf{a}_m = \arg \min_{\mathbf{a}, \beta} \sum_{i=1}^N [\tilde{y}_i - \beta h(\mathbf{x}_i; \mathbf{a})]^2$
5.  $\rho_m = \arg \min_{\rho} \sum_{i=1}^N L(y_i, F_{m-1}(\mathbf{x}_i) + \rho h(\mathbf{x}_i; \mathbf{a}_m))$
6.  $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \rho_m h(\mathbf{x}; \mathbf{a}_m)$

## 7. endFor end Algorithm

感觉该伪代码针对的还是任意损失函数 $L(y_i, F(\mathbf{x}_i))$ 下的回归问题。Algorithm 1 中第 3 步和第 4 步意思是用 $\beta h(\mathbf{x}_i, \mathbf{a})$ 拟合 $F_{m-1}(\mathbf{x})$ 处负梯度，但第 4 步表示只求参数 $\mathbf{a}_m$ ，第 5 步单独求解参数 $\rho_m$ ，这里的疑问是为什么第 4 步要用最小二乘法（即 3.2 节的线性回归）去拟合负梯度（又称伪残差）？

简单理解如下：第 4 步要解的 $h(\mathbf{x}_i, \mathbf{a})$ 相当于梯度下降法中的待解的下降方向 $\mathbf{d}$ ，在梯度下降法中也已提到不必严格限制 $\|\mathbf{d}\|_2 = 1$ ，长度可以由步长 $\alpha$ 调节（例如前面梯度下降法解释中的例 1，若直接取 $d_k = -f'(x_k) = -4$ ，则可得 $\alpha_k = 0.5$ ，仍有 $\Delta x = \alpha_k d_k = -2$ ），因此第 4 步直接用 $h(\mathbf{x}_i, \mathbf{a})$ 拟合负梯度，与梯度下降中约束 $\|\mathbf{d}\|_2 = 1$ 的区别在于未对负梯度除以其模值进行归一化而已。

那为什么不是直接令 $h(\mathbf{x}_i, \mathbf{a})$ 等于负梯度呢？因为这里实际是求假设函数 $h$ ，将数据集中所有的 $\mathbf{x}_i$ 经假设函数 $h$ 映射到对应的伪残差（负梯度） $\tilde{y}_i$ ，所以只能做线性回归了。

李航《统计学习方法》第 8.4.3 节中的算法 8.4 并未显式体现参数 $\rho_m$ ，这应该是第 2 步的(c)步完成的，因为(b)步只是拟合一棵回归树（相当于 Algorithm 1 第 4 步解得 $h(\mathbf{x}_i, \mathbf{a})$ ），而(c)步才确定每个叶结点的取值（相当于 Algorithm 1 第 5 步解得 $\rho_m$ ，只是每个叶结点均对应一个 $\rho_m$ ）；而且回归问题中基函数为实值函数，可以将参数 $\rho_m$ 吸收到基函数中。

## 4、梯度提升树(GBDT)

本部分无实质 GBDT 内容，仅为梳理 GBDT 的概念，具体可参考给出的资源链接。

对于 GBDT，一般资料是按 Gradient Boosting+[CART](#) 处理回归问题讲解的，如林轩田《机器学习技法》课程[第 11 讲](#)（损失函数为平方损失，PPT 中的 C&RT 是 [CART](#) 的基础版本）：

### Gradient Boosted Decision Tree (GBDT)

$s_1 = s_2 = \dots = s_N = 0$

for  $t = 1, 2, \dots, T$

① obtain  $g_t$  by  $\mathcal{A}(\{(\mathbf{x}_n, y_n - s_n)\})$  where  $\mathcal{A}$  is a (squared-error) regression algorithm

—how about sampled and pruned C&RT?

② compute  $\alpha_t = \text{OneVarLinearRegression}(\{(g_t(\mathbf{x}_n), y_n - s_n)\})$

③ update  $s_n \leftarrow s_n + \alpha_t g_t(\mathbf{x}_n)$

return  $G(\mathbf{x}) = \sum_{t=1}^T \alpha_t g_t(\mathbf{x})$

但是，分类问题也可以用回归来处理，例如 3.3 节的对数几率回归，只需将平方损失换为对率损失（参见式(3.27)和式(6.33)，二者关系可参见第 3 章注解中有关式(3.27)的推导）即可。博客[梯度提升树\(GBDT\)原理小结](#)(<https://www.cnblogs.com/pinard/p/6140514.html>)从更一般的角度推导了 GBDT，包括 GBDT 回归版本和分类版本。

特别地，[林轩田](#)老师（个人主页：<https://www.csie.ntu.edu.tw/~htlin/>）的《机器学习基石》和《机器学习技法》两门课程的主页参见 <https://www.csie.ntu.edu.tw/~htlin/mooc/>，课程视频 <https://www.bilibili.com/video/av12463015/>和 <https://www.bilibili.com/video/av12469267/>。

## 5、XGBoost

本部分无实质 XGBoost 内容，仅为梳理 XGBoost 的概念，具体可参考给出的资源链接。

首先，XGBoost 是 **eXtreme Gradient Boosting** 的简称。

其次, XGBoost 与 GBDT 的关系, 可大致类比为 [LIBSVM](#) 与 SVM (或 SMO 算法) 的关系。LIBSVM 是 SVM 算法的一种高效实现软件包, XGBoost 是 GBDT 的一种高效实现; 在实现层面, LIBSVM 对 SMO 算法进行了许多改进, XGBoost 也对 GBDT 进行了许多改进; 另外, LIBSVM 扩展了许多 SVM 变体, XGBoost 也不再仅仅是标准的 GBDT, 也扩展了一些其它功能。

最后, XGBoost 是由[陈天奇](#) (个人主页: <https://homes.cs.washington.edu/~tqchen/>) 开发的; XGBoost [论文](#)发表在 KDD'16 上, 可在 arXiv 下载 <https://arxiv.org/abs/1603.02754>; XGBoost 工具包主页 <https://xgboost.ai/>, 文档 <https://xgboost.readthedocs.io/en/latest/>, 源码 <https://github.com/dmlc/xgboost>; [统计之都](#)(<https://cosx.org/>) 曾对[陈天奇](#)做过一个采访, 参见《[COS 访谈第 18 期: 陈天奇](#)》(<https://cosx.org/2015/06/interview-of-tianqi>), 里面介绍的一些有关作者的信息是论文中看不到的。

官方文档开篇如是介绍 XGBoost:

XGBoost is an optimized distributed gradient boosting library designed to be highly **efficient, flexible and portable**. It implements machine learning algorithms under the **Gradient Boosting** framework. XGBoost provides a parallel tree boosting (also known as GBDT, GBM) that solve many data science problems in a fast and accurate way. The same code runs on major distributed environment (Hadoop, SGE, MPI) and can solve problems beyond billions of examples.

GB、GBDT、XGBoost 的递进关系简单介绍可参见: [一步一步理解 GB、GBDT、xgboost](#) (<https://www.cnblogs.com/wxquare/p/5541414.html>); 另外, 有关 XGBoost 介绍, 很多人推荐[陈天奇的 PPT](#): <https://homes.cs.washington.edu/~tqchen/pdf/BoostedTree.pdf>。

特别崇拜一下 XGBoost 作者陈天奇, 引用新浪微博 [Maybe Shewill](#) 的评价: “撸得了底层框架又搞得了上层高端算法的大神”。

## 8.7 本章小节

本章内容整体上看并不复杂, 可以看成是作者的英文专著《[Ensemble Methods: Foundations and Algorithms](#)》(<https://book.douban.com/subject/10494228/>, 据悉该书中文版即将由电子工业出版社出版) 的精简版, 但 8.2 节有关 AdaBoost 的推导感觉很难吃透, 尤其是本节还涉及到如 Gradient Boosting、GBDT、XGBoost 等概念, 因此必须读一些原始文献和其它参考资料方能对此有所理解, 下面对这些概念进行一个简单的梳理。

**提升(Boosting)**是一族可将弱学习提升为强学习器的算法, 弱学习器间存在强依赖关系、必须串行生成。这族算法的工作机制类似: 先从初始训练集训练出一个基学习器, 再根据基学习器的表现对训练样本分布进行调整, 使得先前基学习器做错的训练样本在后续受到更多关注, 然后基于调整后打着本分布来训练下一个基学习器; 如此重复进行, 直至基学习器数目达到事先指定的值  $T$ , 最终将这  $T$  个基学习器进行加权结合。(摘自西瓜书 P173)

**提升树(Boosting Tree)**是以决策树为基学习器的提升方法(i.e., Boosting + Decision Tree)。对分类问题决策树是二叉分类树, 对回归问题决策树是二叉回归树。(李航《统计学习方法》)

**梯度提升(Gradient Boosting)**是基于加法模型和前向分步算法一种理论框架, 并不限定损失函数, 也不限定基分类器的类型 (类型包含两层含义: 第 1 层如分类和回归等; 第 2 层如决策树、SVM、神经网络等)。

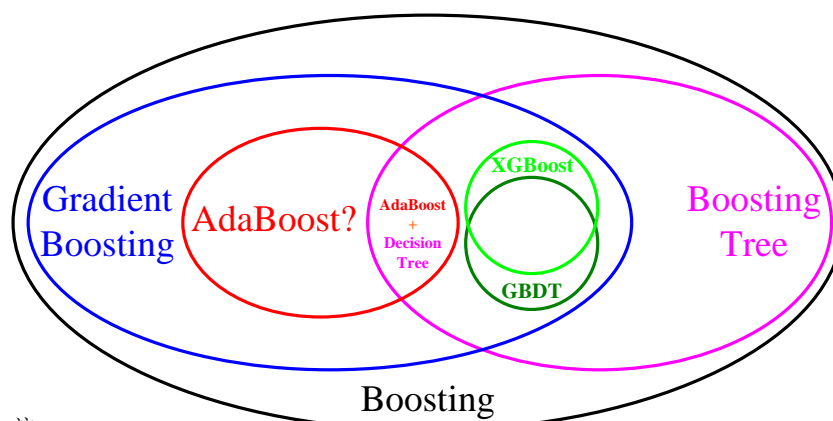
**梯度提升树(GBDT)** = Gradient Boosting + Decision Tree (**CART**), 这里 **CART** 取为回归树。李航《统计学习方法》第 8.4.3 节中的算法 8.4 实际为 GBDT 算法, 与[刘建平 Pinard](#)博客《[梯度提升树\(GBDT\)原理小结](#)》介绍基本一致, 林轩田《[机器学习技法](#)》课程第 11 讲则特指损失函数为平方损失。

特别地，8.2 节介绍的 **AdaBoost(Adaptive Boosting)** 是当 Gradient Boosting 中损失函数为指数损失函数、基分类器为二分类（只能取+1 和-1 两个值）时的特殊形式(i.e., Gradient Boosting + Exponential Loss + Binary base classifier)，这里基分类器可以为任何二分类器，如决策树、SVM、神经网络等。当然，西瓜书作者在第 190 页的阅读材料中提到，这仅是源“统计视角”的一种推导而已，但这派理论产生的推论并不能解释某些现象，因此这可能仅是一个与 AdaBoost 相似的学习过程而非 AdaBoost 本身。

当 AdaBoost 的基分类器为二分类决策树时，也可以称为**提升树**，“对二分类问题，提升树算法只需将 AdaBoost 算法 8.1 中的基本分类器限制为二分类树即可，可以说这时的提升树是 AdaBoost 算法的特殊情况”。但这里与 GBDT 不一样，因为 AdaBoost 限定了基分类器为二分类（实际上也可以做回归任务），并不像 GBDT 里那样，每轮迭代与残差（负梯度）拟合；而且已经提到，GBDT 一般特指 Gradient Boosting 与 CART 回归树的结合。

**XGBoost** 是 GBDT 的一个软件包，类似于 LIBSVM 是 SVM 的一个软件包。XGBoost 在实现 GBDT 时进行了一些细节上的改进和扩展。

各种概念之间的包含关系如下图所示：



注：

- (1) 以上范围重叠仅表示两个概念有交集，但重叠的面积大小与概念的重叠程度不成比例；
- (2) 西瓜书第190页的阅读材料提到8.2节的推导可能仅是一个与AdaBoost相似的学习过程而非AdaBoost本身，因此图中添加了问号；
- (3) XGBoost是GBDT的一种实现，但也有改进和扩展，因此图中二者范围并未完全重叠；

以上仅为个人的初步理解，可能有误，勿被误导^\_^