



## Chapter 8

$$8.1 \quad P(H(n) \leq k) = \sum_{i=0}^k \binom{n}{i} p^i (1-p)^{n-i}$$

Hoeffding inequality:

$$P(H(n) \leq (p-\delta)n) \leq e^{-2\delta n} \quad (8.44)$$

(another version): if  $P(x_i \in [a_i, b_i]) = 1$ ,  $\bar{x} = \frac{x_1 + \dots + x_n}{n}$ ,

$$\text{then } P(\bar{x} - E(\bar{x}) \geq t) \leq \exp \left\{ \frac{-2t^2 n^2}{\sum_{i=1}^n (b_i - a_i)^2} \right\}$$

$$P(|\bar{x} - E(\bar{x})| \geq t) \leq 2 \exp \left\{ \frac{-2t^2 n^2}{\sum_{i=1}^n (b_i - a_i)^2} \right\}.$$

① through (8.44),

$$P(H(x) \neq f(x)) = \sum_{k=0}^{[T/2]} \binom{T}{k} (1-\varepsilon)^k \varepsilon^{(T-k)}$$

$$P(h_i(x) \neq f(x)) = \varepsilon$$

$$H(x) = \text{sign} \left( \sum_{i=1}^T h_i(x) \right)$$

assume that  $k^+(x) \triangleq \sum_{i=1}^T \max(h_i(x), 0)$ ,  $k^-(x) = T - k^+(x)$

$$P(H(x) \neq f(x)) = P(k^+(x) \geq \frac{T}{2})$$

$$= P(k^-(x) < \frac{T}{2}) \quad \text{set } \frac{T}{2} = (\varepsilon - \delta) \cdot T, \text{ that is } \delta = \varepsilon - \frac{1}{2}$$

$$\leq \exp \left\{ -2(\frac{1}{2} - \varepsilon)^2 T \right\}$$

$$= \exp \left\{ -\frac{1}{2} (1-2\varepsilon)^2 T \right\}.$$

② through another expression,

$$x_i \triangleq h_i(x), \quad \bar{x} \geq 0 \text{ denotes that } H(x) = f(x)$$

$$P(H(x) \neq f(x)) = P(\bar{x} < 0)$$

$$E(\bar{x}) = E(h_i(x)) = 1-2\varepsilon$$

$$= P(\bar{x} - E(\bar{x}) < 2\varepsilon - 1)$$

$$\leq \exp \left\{ \frac{-2(2\varepsilon - 1)^2 \cdot T^2}{4T} \right\}$$

$$= \exp \left\{ -\frac{1}{2} (1-2\varepsilon)^2 \cdot T \right\}. \quad \square$$

8.2 没看懂  $L(y_i, F(x))$  是什么东西.

#### 8.4 Gradient Boosting:

$L$  is loss function, for example:  $L_{MSE}(y, F(x)) = \frac{1}{n} \cdot \sum_{i=1}^n (y_i - F(x_i))^2$

$F^*$  is the estimation function,

$$\hat{F}^* = \operatorname{argmin}_{F^*} L(y, F(x))$$

use steepest method, we have

$$\hat{F}_m = \hat{F}_{m-1} - \gamma_m \cdot \left. \frac{\partial L(y, F)}{\partial F} \right|_{F=\hat{F}_{m-1}}$$

we need to calculate  $\hat{h}_m = \left. \frac{\partial L(y, F)}{\partial F} \right|_{F=\hat{F}_{m-1}}$ ,  $\hat{h}_m$  is the basic learner.

1. initializing:  $\hat{F}_0(x) = \operatorname{argmin}_F L(y_i, F)$   $\forall i \in n$

$$2. \gamma_m = \left. \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right|_{F=\hat{F}_{m-1}(x_i)}, \forall i \in n.$$

use  $\{(x_i, \gamma_m)\}_{i=1}^n$  as data, to train  $\hat{h}_m$ .

(Remark: if  $\hat{h}_m(x_i) = \gamma_m$ ,  $\forall i \in n$ .  $\hat{h}_m(x)$  is the best learner)

$$3. \gamma_m = \operatorname{argmin}_{\gamma} L(y, \hat{F}_{m-1}(x) + \gamma \hat{h}_m(x))$$

$$4. \hat{F}_m(x) = \hat{F}_{m-1}(x) + \gamma_m \cdot \hat{h}_m(x).$$

Diff between Gradient Boosting and AdaBoost:

1. Gradient Boost 生成数据  $\{(x_i, \gamma_m)\}$  通过 steepest method and loss function.

得到  $\hat{h}_m$ , 作为下一个 basic learner, 和  $\gamma_m$  参数

2. AdaBoost, 通过 loss function 得到下一个数据集的 weight.

得到 basic function  $\hat{h}_m$  和参数  $\gamma_m$ .

不同点是, Gradient Boosting 类似于 AdaBoost 的推广版本,

可以选择不同的梯度算法, 以及可以适用于 Data 不支持 weight 的情况.  
(Steepest gradient, Newton method)

8.6 因为 Naive Bayesian Classifier 是通过

$$P(x_i|c) \approx \frac{N(x_i|c)}{N(c)} \text{ 出现频率估计的,}$$

Bagging 方法不能降低  $\frac{N(x_i|c)}{N(c)}$  的方差, 所以效果不好  
(方差保持不变)

$$\left\{ \begin{array}{l} \text{Var}\left(\frac{\sum I(x_i|c)}{n}\right) = \frac{1}{n} \text{Var}(I(x_i|c)) \text{ is sustained stably.} \\ P(x_i|c) = E(I(x_i|c)). \end{array} \right.$$

### 8.8. MultiBoosting 算法 和 Iterative Bagging 算法:

1. MultiBoosting : AdaBoost 作为基学习器易过拟合, (over fitted)

Bagging 选取的 Data 需要尽可能平均,

2. Iterative Bagging : Bagging 作为基学习器, 因为 Data 为 weighted data,

当 Data 样本量较小时, 应当停止 Bagging 迭代,

当样本量较小时, 容易欠拟合, (under fitted)