**Lemma 2.3** If $x_1, \ldots, x_m$ are drawn i.i.d. according to $D$ then for any $h$

$$\mathbb{E}_{S \sim D^m} [\hat{R}(h)] = R(h)$$

**Proof** 
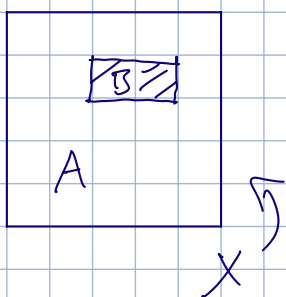$$\mathbb{E}_{S \sim D^m} [\hat{R}(h)] = \mathbb{E}_{S \sim D^m} \frac{1}{m} \sum_{i=1}^{m} [\mathbb{1}_{h(x_i) \neq f(x_i)}]$$

$$= \frac{1}{m} \sum_{i=1}^{m} \mathbb{E} [\mathbb{1}_{h(x_i) \neq f(x_i)}]$$

$$\underset{\text{i.i.d}}{\nearrow} = \frac{1}{m} \sum_{i=1}^{m} \mathbb{E}_{X \sim D} [\mathbb{1}_{h(x) \neq f(x)}] \qquad = \frac{1}{m} \cdot m \, R(h)$$

$$= R(h).$$

## Overfitting

The choice of a suitable hypothesis set $\mathcal{H}$ is important.

**Example** Let $X \subseteq \mathbb{R}^2$ be an axis aligned rectangle, e.g. $X = [0, 1]^2$, let $B \subsetneq X$ be another axis aligned rectangle, let $A = X \setminus B$.



Let $\mathbb{P}$ be a continuous probability distribution on $X$ i.e. $\mathbb{P}(\{x\}) = 0 \ \forall x \in X$. with $\mathbb{P}(A) = \mathbb{P}(B) = \frac{1}{2}$.

$\lceil$ E.g. $X = [0, 1]^2$, $B = [0, \frac{1}{2}] \times [0, 1]$

$$\mathbb{P}(M) = \text{vol}(M) \rfloor$$

Let $f : X \to \{0, 1\}$ be given by

$$f(x) = \begin{cases} 0 & , \ x \in A \\ 1 & , \ x \in B \end{cases}$$

Given a sample $\delta = (x_1, \ldots, x_m)$ with labels $y_i = f(x_i)$ choose hypothesis

$$h_\delta(x) = \begin{cases} y_i & \text{if } x = x_i \text{ for some } i \in [m] \\ 0 & \text{otherwise} \end{cases}$$

Then the empirical risk is

$$\hat{R}(h_S) = \frac{1}{m}\sum_{i=1}^{m} \mathbb{1}_{h_S(x_i) \neq f(x_i)} = 0,$$

since $h_S(x_i) = f(x_i)$ $\forall i \in [m] = \{1, \ldots, m\}$

But the true risk is:

$$\mathbb{P}_{x \sim 1}(h_S(x) \neq f(x)) = \mathbb{P}(x \neq x_1, \ldots, x_m \text{ and } f(x) = 1\})$$

$$= \mathbb{P}(B \setminus \{x_1, \ldots, x_m\})$$

$$\underset{\uparrow}{=} \mathbb{P}(B) = \frac{1}{2}$$

$$\mathbb{P}(\{x_1, \ldots, x_m\}) = 0 \text{ since } \mathbb{P}(\{x\}) = 0 \ \forall x \in X$$

Hence $\hat{R}(h_S)$ is minimal $(=0)$, but predictor

is bad $R(h_S) = \mathbb{P}(h_S(x) \neq f(x)) = \frac{1}{2}$

not better than random guessing.

# Thm 2.6

Let $H$ be a finite set of functions $h: X \rightarrow \{0,1\}$
Assume $f \in H$ and let $A$ be an algorithm
that for each i.i.d sample $S = (x_1, \ldots, x_m)$
and labeled training data $(x_i, f(x_i))$, $i = 1, \ldots, m$
returns an $h_S \in H$ with $\hat{R}(h_S) = 0$.
Then for $\varepsilon, \delta > 0$ if

$$m \geq \frac{1}{\varepsilon}\left(\log |H| + \log\left(\frac{1}{\delta}\right)\right)$$

then with prob. $\geq 1-\delta$

$$R(h_S) \leq \varepsilon.$$

## Proof

$R(h_S)$ depends on training sample $S$ and is
difficult to evaluate directly. Instead, we bound it
as follows:

Let $0 < \varepsilon < 1$, and let $H = \{h_1, \ldots, h_n\}$
$$n = |H| = \#H$$

$$\mathbb{P}(R(h_S) > \varepsilon) = \mathbb{P}(\hat{R}(h_S) = 0 \text{ and } R(h_S) > \varepsilon)$$

$$\leq \mathbb{P}(\hat{R}(h) = 0 \text{ and } R(h) > \varepsilon \quad \text{for some } h \in H)$$

$$= \mathbb{P}\left(\{\hat{R}(h_1) = 0 \ \& \ R(h_1) > \varepsilon\} \cup \{\hat{R}(h_2) = 0 \ \& \ R(h_2) > \varepsilon\}\right.$$

$$\left. \cup \cdots \cup \{\hat{R}(h_n) = 0 \ \& \ R(h_n) > \varepsilon\}\right)$$

$$\leq \sum_{h \in \mathcal{H}} \mathbb{P}\left(\{\hat{R}(h) = 0\} \cap \underbrace{\{R(h) > \varepsilon\}}_{\text{deterministic event (indep. of } S)}\right)$$

$$= \sum_{\substack{h \in \mathcal{H} \\ R(h) > \varepsilon}} \mathbb{P}\left(\hat{R}(h) = 0\right)$$

Fix $h \in \mathcal{H}$ with $R(h) > \varepsilon$ $\qquad$ and recall

$$\hat{R}(h) = \frac{1}{m} \sum_{j=1}^{m} \mathbb{1}_{h(x_j) \neq f(x_j)}$$

$$\Rightarrow \mathbb{P}\left(\hat{R}(h) = 0\right) = \mathbb{P}\left(h(x_j) = f(x_j) \text{ for } j = 1, \ldots, m\right)$$

$$\overset{i.i.d.}{=} \prod_{j=1}^{m} \mathbb{P}\left(h(x_j) = f(x_j)\right)$$

$$= \prod_{j=1}^{m} \left(1 - \underbrace{\mathbb{P}\left(h(x_j) \neq f(x_j)\right)}_{= R(h) > \varepsilon}\right)$$

$$< \prod_{j=1}^{m} (1 - \varepsilon) = (1 - \varepsilon)^m$$

$$\Rightarrow \quad \mathbb{P}\left(R(h_S) > \varepsilon\right) \quad \leq \quad \sum_{\substack{h \in \mathcal{H} \\ R(h) > \varepsilon}} (1-\varepsilon)^m \quad \leq \quad |\mathcal{H}|\,(1-\varepsilon)^m$$

$$\leq \quad |\mathcal{H}|\,e^{-\varepsilon m}$$

$$\uparrow$$

$$1-x \leq e^{-x} \quad \forall x \in \mathbb{R}$$

$$\rightsquigarrow \quad \mathbb{P}\left(R(h_S) \leq \varepsilon\right) \quad \geq \quad 1 - |\mathcal{H}|\,e^{-\varepsilon m}$$

With $\quad \delta := |\mathcal{H}|\,e^{-\varepsilon m}\quad$ we obtain

$$\mathbb{P}\left(R(h_S) \leq \varepsilon\right) \quad \geq \quad 1-\delta$$

$$\Rightarrow$$

For $\delta \in (0,1)$ and $\quad m \geq \frac{1}{\varepsilon}\log\left(|\mathcal{H}|/\delta\right) = \frac{1}{\varepsilon}\left(\log|\mathcal{H}| + \log\left(\tfrac{1}{\delta}\right)\right)$

$$\mathbb{P}\left(R(h_S) \leq \varepsilon\right) \quad \geq \quad 1-\delta$$

$\quad$ solve $\quad \delta = |\mathcal{H}|\,e^{-\varepsilon m}\quad$ for $m$

$$\Rightarrow \quad e^{\varepsilon m} = \frac{|\mathcal{H}|}{\delta}$$

$$\Rightarrow \quad m = \frac{1}{\varepsilon}\log\left(|\mathcal{H}|/\delta\right)$$

**Conclusion** For a finite hypothesis set $\mathcal{H}$, a consistent learning algorithm $A$ is a PAC-learning algorithm with sample complexity polynomial (even linear) in $1/\varepsilon$ and logarithmic in $|\mathcal{H}|$ and $1/\delta$.

$\log|\mathcal{H}|$ may be interpreted as the number of bits required to represent $\mathcal{H}$ (up to a constant factor).

Note that $f \in \mathcal{H}$ guarantees that ERM always returns an $h_S$ with $\hat{R}(h_S) = 0$.