

These notes use the notation, definitions, and results from:

- *Bias of a random coin*[BRC] <http://tinyurl.com/bias-rand-coin> ¹
- *Pinsker's inequality*[Pinsker] <http://tinyurl.com/pol-pin> ²

For any distribution D over $\mathcal{X} \times \{-1, 1\}$ and hypothesis class \mathcal{H} , define

$$h_D^* := \arg \min_{h \in \mathcal{H}} \text{err}(h, D)$$

(there is no loss of generality in assuming that the minimum is achieved and unique).

For any $h \in \mathcal{H}$, define its **excess risk** under D :

$$\Delta(h, D) := \text{err}(h, D) - \text{err}(h_D^*, D).$$

Note that $\Delta(h, D) \geq 0$ always (why?).

Theorem 0.1 *For any hypothesis class \mathcal{H} with VC-dimension $d \geq 1$ and $n \geq d/\log 3$, we have*

$$\inf_{\hat{h}_n} \sup_D \mathbb{E}_{D^n} [\Delta(\hat{h}_n, D)] \geq \sqrt{\frac{d/n}{64 \log 3}} \approx 0.1193 \sqrt{\frac{d}{n}},$$

where the infimum is over all learning algorithms \mathcal{L} that map samples of size n to hypotheses $\hat{h}_n \in \mathcal{H}$ and the supremum is over all distributions over the instance space \mathcal{X} .

The quantity on the left-hand side of the inequality is known as the minimax expected excess risk.

Proof:

Step I: Defining the adversarial distribution. There is no loss of generality in taking $\mathcal{X} = [d] \equiv \{1, \dots, d\}$ and $\mathcal{H} = \{-1, 1\}^{\mathcal{X}}$ (why?). For $\gamma \in \{-1, 1\}^{\mathcal{X}}$ and $b \in [0, 1]$, define the distribution $D_{b,\gamma}$ over $\mathcal{X} \times \{-1, 1\}$ as follows:

$$\mathbb{P}_{(X,Y) \sim D_{b,\gamma}} (X = x, Y = y) = \frac{1}{d} \left(\frac{1}{2} + \frac{yb\gamma(x)}{2} \right).$$

In words, X is drawn uniformly from $\mathcal{X} = [d]$ and conditioned on $X = x$, Y is drawn from $\text{R}(b\gamma(x))$, where $\text{R}(\cdot)$ is the Rademacher distribution defined in [BRC].

¹long url: <https://www.cs.bgu.ac.il/~asml162/wiki.files/b-k-opt.pdf>

²long url: <https://www.cs.bgu.ac.il/~asml162/wiki.files/pollard-pinsker.pdf>

Step II: Computing the excess risk. Convince yourself that for $D_{g,\gamma}$ as defined above,

$$h_{D_{b,\gamma}}^*(x) = \arg \max_{y \in \{-1,1\}} \mathbb{P}_{(X,Y) \sim D_{b,\gamma}} (Y = y | X = x) = \gamma(x)$$

(this follows directly from (1) and (2) in [BRC]). Since $h_{D_{b,\gamma}}^*$ does not depend on b , we will write $h_\gamma^* := h_{D_{b,\gamma}}^*$ as a shorthand. Let us compute $\text{err}(h_\gamma^*, D_{b,\gamma})$:

$$\begin{aligned} \text{err}(h_\gamma^*, D_{b,\gamma}) &= \mathbb{P}_{(X,Y) \sim D_{b,\gamma}} (\gamma(X) \neq Y) \\ &= \mathbb{P}(\gamma(X) = 1) \mathbb{P}(Y \neq 1 | \gamma(X) = 1) + \mathbb{P}(\gamma(X) = -1) \mathbb{P}(Y \neq -1 | \gamma(X) = -1) \\ &= \frac{1}{2} - \frac{b}{2} \quad (\text{ex1}) \end{aligned}$$

(exercise: verify this carefully).

Now let us compute the $\text{err}(h, D_{b,\gamma})$ for an arbitrary $h \in \mathcal{H}$:

$$\begin{aligned} \text{err}(h, D_{b,\gamma}) &= \mathbb{P}_{(X,Y) \sim D_{b,\gamma}} (h(X) \neq Y) \\ &= \mathbb{P}(h(X) = h_\gamma^*(X)) \cdot \mathbb{P}(h(X) \neq Y | h(X) = h_\gamma^*(X)) \\ &\quad + \mathbb{P}(h(X) \neq h_\gamma^*(X)) \cdot \mathbb{P}(h(X) \neq Y | h(X) \neq h_\gamma^*(X)). \end{aligned}$$

We know from the above that

$$\begin{aligned} \mathbb{P}(h(X) \neq Y | h(X) = h_\gamma^*(X)) &= \mathbb{P}(h_\gamma^*(X) \neq Y) = \frac{1}{2} - \frac{b}{2}, \\ \mathbb{P}(h(X) \neq Y | h(X) \neq h_\gamma^*(X)) &= \mathbb{P}(h_\gamma^*(X) = Y) = \frac{1}{2} + \frac{b}{2}, \end{aligned}$$

whence

$$\begin{aligned} \text{err}(h, D_{b,\gamma}) &= \mathbb{P}(h(X) = h_\gamma^*(X)) \left(\frac{1}{2} - \frac{b}{2} \right) + \mathbb{P}(h(X) \neq h_\gamma^*(X)) \left(\frac{1}{2} + \frac{b}{2} \right) \\ &= \left(\frac{1}{2} - \frac{b}{2} \right) + b \mathbb{P}(h(X) \neq h_\gamma^*(X)) \\ &= \text{err}(h_\gamma^*, D_{b,\gamma}) + \frac{b}{d} \sum_{x \in [d]} \mathbb{1}_{\{h(x) \neq \gamma(x)\}} \quad (\text{ex2}) \end{aligned}$$

(exercise: verify this carefully).

We conclude that

$$\Delta(h, D_{b,\gamma}) = \frac{b}{d} \sum_{x \in [d]} \mathbb{1}_{\{h(x) \neq \gamma(x)\}}. \quad (1)$$

Step III: Taking expectations. Let $S_n(X_i, Y_i)_{i \in [n]}$ be drawn from $D_{b,\gamma}^n$ and define the random variables N_x to be the number of occurrences of the point $x \in \mathcal{X}$ in the random sample S_n :

$$N_x = \sum_{i=1}^n \mathbb{1}_{\{X_i=x\}}. \quad (2)$$

Let $\mathcal{L} : S_n \mapsto \hat{h}_n \in \mathcal{H}$ be any learning algorithm, and let us lower-bound $\mathbb{E}[\Delta(\hat{h}_n, D_{b,\gamma})]$, where the expectation is both over $\gamma \sim \text{Unif}(\{-1, 1\}^d)$ and $S_n \sim D_{b,\gamma}^n$:

$$\begin{aligned} \mathbb{E}_{\gamma \sim \text{Unif}(\{-1, 1\}^d)} \mathbb{E}_{S_n \sim D_{b,\gamma}^n} [\Delta(\hat{h}_n, D_{b,\gamma})] &= \mathbb{E}_{\gamma} \mathbb{E}_{S_n} \left[\frac{b}{d} \sum_{x \in [d]} \mathbb{1}_{\{\hat{h}_n(x) \neq \gamma(x)\}} \right] \\ &= \frac{b}{d} \sum_{x \in [d]} \mathbb{P}_{\gamma \sim \text{Unif}(\{-1, 1\}^d), \hat{h}_n \sim \mathcal{L}(S_n)} (\hat{h}_n(x) \neq \gamma(x)) \end{aligned}$$

— note that the inner $\mathbb{P}(\cdot)$ is NOT over the input to \hat{h}_n but rather over \hat{h}_n itself! Let us focus on a fixed $x \in [d]$:

$$\begin{aligned} \mathbb{P}_{\gamma \sim \text{Unif}(\{-1, 1\}^d), \hat{h}_n \sim \mathcal{L}(S_n)} (\hat{h}_n(x) \neq \gamma(x)) &= \mathbb{E}_{N_x} \left[\mathbb{P}_{\gamma, \hat{h}_n} (\hat{h}_n(x) \neq \gamma(x)) \mid N_x \right] \\ &= \sum_{k=0}^m \mathbb{P}(N_x = k) \mathbb{P}_{\gamma, \hat{h}_n} (\hat{h}_n(x) \neq \gamma(x) \mid N_x = k), \end{aligned}$$

where we can pull the \mathbb{E}_{γ} inside because N_x and γ are independent. Recall the definition of $\text{opt}(k, b)$ from [BRC] and convince yourself that

$$\mathbb{P}_{\gamma \sim \text{Unif}(\{-1, 1\}^d), \hat{h}_n \sim \mathcal{L}(S_n)} (\hat{h}_n(x) \neq \gamma(x) \mid N_x = k) \geq \text{opt}(k, b); \quad (\text{ex3})$$

this is the most important step of the proof.

From here, it is easy to conclude that

$$\mathbb{E}_{\gamma \sim \text{Unif}(\{-1, 1\}^d)} \mathbb{E}_{S_n \sim D_{b,\gamma}^n} [\Delta(\hat{h}_n, D_{b,\gamma})] \geq \frac{b}{d} \sum_{x \in [d]} \mathbb{E}_{N_x} \text{opt}(N_x, b). \quad (3)$$

Step IV: Lower-bounding opt. Recall from [BRC] that

$$\text{opt}(k, b) = \frac{1}{2} \left(1 - \frac{1}{2} \left\| \mathbb{R}(b)^k - \mathbb{R}(-b)^k \right\|_1 \right) \quad (4)$$

and from [Pinsker] that for all distributions P, Q ,

$$\|P - Q\|_1 \leq \sqrt{2 \text{KL}(P||Q)}. \quad (5)$$

We will apply (5) to $P = \mathbb{R}(b)^k$, $Q = \mathbb{R}(-b)^k$. From Homework 3.2(b), we know that

$$\begin{aligned} \text{KL}(\mathbb{R}(b)^k || \mathbb{R}(-b)^k) &= k \text{KL}(\mathbb{R}(b) || \mathbb{R}(-b)) \\ &= kb \log \frac{1+b}{1-b} \end{aligned} \quad (\text{ex4})$$

(exercise: verify the last claim). From Homework 3.3, we know that

$$\sqrt{2b \log \frac{1+b}{1-b}} \leq 2b\sqrt{\log 3}, \quad 0 \leq b \leq \frac{1}{2} \quad (6)$$

and hence

$$\text{opt}(k, b) \geq \frac{1}{2} \left(1 - \sqrt{(\log 3)kb^2} \right). \quad (7)$$

Substituting into (3), we get

$$\mathbb{E}_{\gamma \sim \text{Unif}(\{-1,1\}^d)} \mathbb{E}_{S_n \sim D_{b,\gamma}^n} \left[\Delta(\hat{h}_n, D_{b,\gamma}) \right] \geq \frac{b}{2d} \sum_{x \in [d]} \mathbb{E}_{N_x} \left[1 - \sqrt{(\log 3)N_x b^2} \right].$$

Now the function $t \mapsto \left(1 - \sqrt{(\log 3)tb^2} \right)$ is convex, and furthermore, $\mathbb{E}[N_x] = n/d$. Hence, by Jensen's inequality,

$$\frac{b}{2d} \sum_{x \in [d]} \mathbb{E}_{N_x} \left(1 - \sqrt{(\log 3)N_x b^2} \right) \geq \frac{b}{2} \left(1 - \sqrt{\frac{(\log 3)nb^2}{d}} \right) =: B_{n,d}(b)$$

(what happened to $\frac{1}{d} \sum_x x$?).

Step V: optimizing b . It remains specify the parameter b , which we choose to be of the form $b = c\sqrt{d/n}$, where c will be chosen so as to maximize $B_{n,d}(\cdot)$. Now

$$B_{n,d} \left(c\sqrt{d/n} \right) = \frac{1}{2} \sqrt{\frac{d}{n}} c \left(1 - c\sqrt{\log 3} \right),$$

which achieves the maximum value of $\sqrt{\frac{d/n}{64 \log 3}}$ at $c = \frac{1}{\sqrt{4 \log 3}}$ [(ex5): verify this!].

This proves the claim, since

$$\sup_D \mathbb{E}_{D^n} \left[\Delta(\hat{h}_n, D) \right] \geq \mathbb{E}_{\gamma} \mathbb{E}_{D_{b,\gamma}^n} \left[\Delta(\hat{h}_n, D_{b,\gamma}) \right].$$

The condition $n \geq d/\log 3$ ensures that $b = \sqrt{\frac{d}{4(\log 3)n}} \leq \frac{1}{2}$, so that the estimate in (6) applies. ■