

Minimax lower bounds

Maxim Raginsky

December 4, 2013

Now that we have a good handle on the performance of ERM and its variants, it is time to ask whether we can do better. For example, consider binary classification: we observe n i.i.d. training samples from an unknown joint distribution P on $X \times \{0, 1\}$, where X is some feature space, and for a fixed class \mathcal{F} of candidate classifiers $f: X \rightarrow \{0, 1\}$ we let \hat{f}_n be the ERM solution

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{f(X_i) \neq Y_i\}}. \quad (1)$$

IF \mathcal{F} is a VC class with VC dimension V , then the excess risk of \hat{f}_n over the best-in-class performance $L^*(\mathcal{F}) \equiv \inf_{f \in \mathcal{F}} L(f)$ satisfies

$$L(\hat{f}_n) \leq L^*(\mathcal{F}) + C \left(\sqrt{\frac{V}{n}} + \sqrt{\frac{\log(1/\delta)}{n}} \right)$$

with probability at least $1 - \delta$, where $C > 0$ is some absolute constant. Integrating, we also get the following bound on the expected excess risk:

$$\mathbb{E} [L(\hat{f}_n) - L^*(\mathcal{F})] \leq C \sqrt{\frac{V}{n}}, \quad (2)$$

for some constant $C > 0$. Crucially, the bound (2) holds for *all* possible joint distributions P on $X \times \{0, 1\}$, and the right-hand side is *independent* of P — it depends only on the properties of the class \mathcal{F} ! Thus, we deduce the following remarkable *distribution-free guarantee* for ERM: for any VC class \mathcal{F} , the ERM algorithm (1) satisfies

$$\sup_{P \in \mathcal{D}(X \times \{0, 1\})} \mathbb{E}_P [L_P(\hat{f}_n) - L_P^*(\mathcal{F})] \leq C \sqrt{\frac{V}{n}}. \quad (3)$$

(We have used subscript P to explicitly indicate the fact that the quantity under the supremum depends on the underlying distribution P . In the sequel, we will often drop the subscript to keep the notation uncluttered.) Let's take a moment to reflect on the significance of the bound (3). What it says is that, regardless of how "weird" the stochastic relationship between the feature $X \in X$ and the label $Y \in \{0, 1\}$ might be, as long as we scale our ambition back and aim at approaching the performance of the best classifier in some VC class \mathcal{F} , the ERM algorithm will produce a good classifier with a uniform $O(\sqrt{V/n})$ guarantee on its excess risk.

At this point, we stop and ask ourselves: could this bound be too pessimistic, even when we are so lucky that the optimal (Bayes) classifier happens to be in \mathcal{F} ? (Recall that the Bayes classifier for a given P has the form

$$f^*(x) = \begin{cases} 1, & \text{if } \eta(x) \geq 1/2 \\ 0, & \text{otherwise} \end{cases},$$

where $\eta(x) = \mathbb{E}[Y|X = x] = \mathbb{P}(Y = 1|X = x)$ is the *regression function*.) Let $\mathcal{D}(\mathcal{F})$ denote the subset of $\mathcal{D}(X \times \{0, 1\})$ consisting of all joint distributions of $X \in X$ and $Y \in \{0, 1\}$, such that $f^* \in \mathcal{F}$. Then from (2) we have

$$\sup_{P \in \mathcal{D}(\mathcal{F})} \mathbb{E}[L(\hat{f}_n) - L(f^*)] \leq C \sqrt{\frac{V}{n}}, \quad (4)$$

where \hat{f}_n is the ERM solution (1). However, we know that if the relationship between X and Y is deterministic, i.e., if $Y = f^*(X)$, then ERM performs much better. More precisely, let $\mathcal{D}_0(\mathcal{F})$ be the *zero-error class*:

$$\mathcal{D}_0(\mathcal{F}) \triangleq \{P \in \mathcal{D}(\mathcal{F}) : Y = f^*(X) \text{ a.s.}\}.$$

Then one can show that

$$\sup_{P \in \mathcal{D}_0(\mathcal{F})} \mathbb{E}[L(\hat{f}_n) - L(f^*)] \leq C \frac{V}{n}, \quad (5)$$

a much better bound than the “global” bound (4) (see, e.g., the book by Vapnik [Vap98]). This suggests that the performance of ERM is somehow tied to how “sharp” the behavior of η is around the decision boundary that separates the sets $\{x \in X : \eta(x) \geq 1/2\}$ and $\{x \in X : \eta(x) < 1/2\}$. To see whether this is the case, let us define, for any $h \in [0, 1]$, the class of distributions

$$\mathcal{D}(h, \mathcal{F}) \triangleq \{P \in \mathcal{D}(\mathcal{F}) : |2\eta(X) - 1| \geq h \text{ a.s.}\}$$

(in that case, the distributions in $\mathcal{D}(h, \mathcal{F})$ are said to satisfy the *Massart noise condition* with margin h .) We have already seen the two extreme cases:

- $h = 0$ — this gives $\mathcal{D}(0, \mathcal{F}) = \mathcal{D}(\mathcal{F})$ (the bound $|2\eta - 1| \geq 0$ holds trivially for any P).
- $h = 1$ — this gives the zero-error regime $\mathcal{D}(1, \mathcal{F}) = \mathcal{D}_0(\mathcal{F})$ (if $|2\eta - 1| \geq 1$ a.s., then η can take only values 0 and 1 a.s.).

However, intermediate values of h are also of interest: if a distribution P belongs to $\mathcal{D}(h, \mathcal{F})$ for some $0 < h < 1$, then its regression function η makes a jump of size h as we cross the decision boundary. With this in mind, for any $n \in \mathbb{N}$ and $h \in [0, 1]$ let us define the *minimax risk*

$$R_n(h, \mathcal{F}) \triangleq \inf_{\tilde{f}_n} \sup_{P \in \mathcal{D}(h, \mathcal{F})} \mathbb{E}[L(\tilde{f}_n) - L(f^*)], \quad (6)$$

where the infimum is over *all* learning algorithms \tilde{f}_n based on n i.i.d. training samples. The term “minimax” indicates that we are *minimizing* over all admissible learning algorithms, while *maximizing* over all distributions in a given class. The following result was proved by Pascal Massart and Élodie Nédélec [MN06]:

Theorem 1. Let \mathcal{F} be a VC class of binary-valued functions on X with VC dimension $V \geq 2$. Then for any $n \geq V$ and any $h \in [0, 1]$ we have the lower bound

$$R_n(h, \mathcal{F}) \geq c \min \left(\sqrt{\frac{V}{n}}, \frac{V}{nh} \right), \quad (7)$$

where $c > 0$ is some absolute constant.

Let us examine some implications:

- When $h = 0$, the right-hand side of (7) is equal to $c\sqrt{V/n}$. Thus, without any further assumptions, ERM is as good as it gets (it is minimax-optimal), up to multiplicative constants.
- When $h = 1$ (the zero-error case), the right-hand side of (7) is equal to cV/n , which matches the upper bound (5) up to constants. Thus, if we happen to know that we are in a zero-error situation, ERM is minimax-optimal as well.
- For intermediate values of h , the lower bound depends on the relative sizes of h , V , and n . In particular, if $h \geq \sqrt{V/n}$, we have the minimax lower bound $R_n(h, \mathcal{F}) \geq cV/nh$. Alternatively, for a fixed $h \in (0, 1)$, we may think of $n^* = \lceil V/h^2 \rceil$ as the *cutoff sample size*, beyond which the effect of the margin condition on η can be “spotted” and exploited by a learning algorithm.
- In the same paper, Massart and Nédélec obtain the following upper bound on ERM:

$$\sup_{P \in \mathcal{P}(h, \mathcal{F})} \mathbb{E} [L(\hat{f}_n) - L(f^*)] \leq \begin{cases} C \sqrt{\frac{V}{n}}, & \text{if } h \leq \sqrt{V/n} \\ C \frac{V}{nh} \left(1 + \log \frac{nh^2}{V} \right), & \text{if } h > \sqrt{V/n} \end{cases}. \quad (8)$$

Thus, ERM is nearly minimax-optimal (we say “nearly” because of the extra log factor in the above bound; in fact, as Massart and Nédélec show, the log factor is unavoidable when the function class \mathcal{F} is “rich” in a certain sense). The proof of the above upper bound is rather involved and technical, and we will not get into it here.

The appearance of the logarithmic term in (8) is rather curious. Given the lower bound of Theorem 1, one may be tempted to dismiss it as an artifact of the analysis used by Massart and Nédélec. However, as we will now see, in certain situations this logarithmic term is unavoidable. To that end, we first need a definition: We say that a class \mathcal{F} of binary-valued functions $f : X \rightarrow \{0, 1\}$ is (N, D) -rich, for some $N, D \in \mathbb{N}$, if there exist N distinct points $x_1, \dots, x_N \in X$, such that the projection

$$\mathcal{F}(x^N) = \left\{ (f(x_1), \dots, f(x_N)) : f \in \mathcal{F} \right\}$$

of \mathcal{F} onto x^N contains all binary strings of Hamming weight¹ D . Some examples:

- If \mathcal{F} is a VC-class with VC dimension V , then it is (V, D) -rich for all $1 \leq D \leq V$. This follows directly from definitions.

¹The Hamming weight of a binary string is the number of nonzero bits it has.

- A nontrivial example, and one that is relevant to statistical learning, is as follows. Let \mathcal{F} be the collection of indicators of all halfspaces in \mathbb{R}^d , for some $d \geq 2$. There is a result in computational geometry which says that, for any $N \geq d + 1$, one can find N distinct points x_1, \dots, x_N , such that $\mathcal{F}(x^n)$ contains all strings in $\{0, 1\}^N$ with Hamming weight up to, and including, $\lfloor d/2 \rfloor$. Consequently, \mathcal{F} is $(N, \lfloor d/2 \rfloor)$ -rich for all $N \geq d + 1$.

We can now state the following result [MN06]:

Theorem 2. *Given some $D \geq 1$, suppose that \mathcal{F} is (N, D) -rich for all $N \geq 4D$. Then*

$$R_n(h, \mathcal{F}) \geq c(1-h) \frac{D}{nh} \left[1 + \log \frac{nh^2}{D} \right] \quad (9)$$

for any $\sqrt{D/n} \leq h < 1$, where $c > 0$ is some absolute constant.

We will present the proofs of Theorems 1 and 2 in Sections 2 and 3, after giving some necessary background on minimax lower bounds.

1 Preparation: minimax lower bounds for statistical estimation

To prove Theorem 1, we need some background on minimax lower bounds for statistical estimation problems. The presentation here follows an excellent paper by Bin Yu [Yu97].

The setting is as follows: we have an indexed set $\{P_\theta : \theta \in \Theta\}$ of probability distributions on some set Z , where Θ is some parameter set. For our purposes, it suffices to consider the case when Z is a finite set. We observe a random sample Z from one of the P_θ 's, and our goal is to estimate θ . We will measure the quality of any estimator $\hat{\theta} = \hat{\theta}(Z)$ by its expected risk

$$\mathbb{E}_\theta[d(\theta, \hat{\theta}(Z))] = \sum_{z \in Z} P_\theta(z) d(\theta, \hat{\theta}(z)),$$

where $\mathbb{E}_\theta[\cdot]$ denotes expectation with respect to P_θ , and $d : \Theta \times \Theta \rightarrow \mathbb{R}^+$ is a given loss function. Here, we will assume that d is a *pseudometric* on Θ , i.e., it has the following properties:

1. Symmetry – $d(\theta, \theta') = d(\theta', \theta)$ for any $\theta, \theta' \in \Theta$.
2. Triangle inequality – $d(\theta, \theta') \leq d(\theta, \eta) + d(\eta, \theta')$ for any $\theta, \theta', \eta \in \Theta$.

(We do not require that $d(\theta, \theta') = 0$ if and only if $\theta = \theta'$, in which case d is a *metric*.) Since we do not know ahead of time which of the θ 's we will be facing, it is natural to expect the worst and seek an estimator $\hat{\theta}$ to minimize the *worst-case* risk $\sup_{\theta \in \Theta} \mathbb{E}_\theta[d(\theta, \hat{\theta}(Z))]$. With this in mind, let us define the *minimax risk*

$$\mathfrak{M}(\Theta) \triangleq \inf_{\hat{\theta}} \sup_{\theta \in \Theta} \mathbb{E}_\theta[d(\theta, \hat{\theta}(Z))],$$

where the infimum is over all estimators $\hat{\theta} = \hat{\theta}(Z)$. We are particularly interested in tight *lower bounds* on $\mathfrak{M}(\Theta)$, since they will give us some idea of how difficult the estimation problem is.

A simple yet powerful technique for getting lower bounds is the *two-point method* introduced by Lucien Le Cam:

Lemma 1. For any two $\theta, \theta' \in \Theta$ and any estimator $\hat{\theta}$,

$$\mathbb{E}_\theta [d(\theta, \hat{\theta}(Z))] + \mathbb{E}_{\theta'} [d(\theta', \hat{\theta}(Z))] \geq d(\theta, \theta') \cdot \sum_{z \in Z} \min(P_\theta(z), P_{\theta'}(z)). \quad (10)$$

Proof. Consider an arbitrary $z \in Z$. If $P_\theta(z) \leq P_{\theta'}(z)$, then

$$\begin{aligned} P_\theta(z)d(\theta, \hat{\theta}(z)) + P_{\theta'}d(\theta', \hat{\theta}(z)) &= P_\theta(z) [d(\theta, \hat{\theta}(z)) + d(\theta', \hat{\theta}(z))] + [P_{\theta'}(z) - P_\theta(z)] d(\theta', \hat{\theta}(z)) \\ &\geq P_\theta(z) [d(\theta, \hat{\theta}(z)) + d(\theta', \hat{\theta}(z))] \\ &\geq P_\theta(z)d(\theta, \theta'), \end{aligned}$$

where the second line is because $P_{\theta'}(z) \geq P_\theta(z)$ and $d(\cdot, \cdot)$ is nonnegative, while the third line is by the triangle inequality. By the same token, if $P_\theta(z) \geq P_{\theta'}(z)$, then

$$P_\theta(z)d(\theta, \hat{\theta}(z)) + P_{\theta'}d(\theta', \hat{\theta}(z)) \geq P_{\theta'}(z)d(\theta, \theta').$$

Summing over all $z \in Z$, we get (10). □

The sum on the right-hand side of (10) can be expressed in terms of the so-called *total variation distance* between P_θ and $P_{\theta'}$. For any two probability distributions $P, Q \in \mathcal{P}(Z)$, the total variation distance is

$$\|P - Q\|_{\text{TV}} \triangleq \frac{1}{2} \sum_{z \in Z} |P(z) - Q(z)|. \quad (11)$$

Moreover,

$$\|P - Q\|_{\text{TV}} = 1 - \sum_{z \in Z} \min(P(z), Q(z)). \quad (12)$$

To prove (12), we just use the definition (11): if we let $Z' \triangleq \{z \in Z : P(z) \geq Q(z)\}$, then

$$\begin{aligned} \|P - Q\|_{\text{TV}} &= \frac{1}{2} \sum_{z \in Z'} (P(z) - Q(z)) + \frac{1}{2} \sum_{z \notin Z'} (Q(z) - P(z)) \\ &= \frac{1}{2} (P(Z') - Q(Z')) + \frac{1}{2} (Q(Z'^c) - P(Z'^c)) \\ &= \frac{1}{2} (P(Z') - Q(Z')) + \frac{1}{2} (P(Z') - Q(Z')) \\ &= P(Z') - Q(Z') \\ &= \sum_{z: P(z) \geq Q(z)} (P(z) - Q(z)) \\ &= \sum_{z: P(z) \geq Q(z)} (P(z) - \min(P(z), Q(z))) \\ &= \sum_{z \in Z} (P(z) - \min(P(z), Q(z))) - \underbrace{\sum_{z: P(z) \leq Q(z)} (P(z) - \min(P(z), Q(z)))}_{=0} \\ &= 1 - \sum_{z \in Z} \min(P(z), Q(z)). \end{aligned}$$

and we are done. Thus, we can rewrite the bound (10) as

$$\mathbb{E}_\theta [d(\theta, \hat{\theta}(Z))] + \mathbb{E}_{\theta'} [d(\theta', \hat{\theta}(Z))] \geq d(\theta, \theta') \cdot (1 - \|P_\theta - P_{\theta'}\|_{\text{TV}}). \quad (13)$$

Let us examine some consequences. If we take the supremum of both sides of (13) over the choices of θ and θ' , then we see that, for any estimator $\hat{\theta}$,

$$\sup_{\theta \in \Theta} \mathbb{E}_{\theta} [d(\theta, \hat{\theta}(Z))] \geq \frac{1}{2} \sup_{\theta \neq \theta'} \{d(\theta, \theta') \cdot (1 - \|P_{\theta} - P_{\theta'}\|_{\text{TV}})\}. \quad (14)$$

Notice that the right-hand side does not depend on the choice of $\hat{\theta}$, so we can take the infimum of both sides of (14) over $\hat{\theta}$ and obtain the following:

Corollary 1. *The minimax risk $\mathfrak{M}(\Theta)$ can be lower-bounded as*

$$\mathfrak{M}(\Theta) \geq \frac{1}{2} \sup_{\theta \neq \theta'} \{d(\theta, \theta') \cdot (1 - \|P_{\theta} - P_{\theta'}\|_{\text{TV}})\}. \quad (15)$$

The “two-point” bound (15) gives us an idea of when the problem of estimating θ based on observations $Z \sim P_{\theta}$ is difficult: when there exists at least one pair of distinct parameter values $\theta, \theta' \in \Theta$ such that $d(\theta, \theta')$ is large, while the total variation distance $\|P_{\theta} - P_{\theta'}\|_{\text{TV}}$ is small. This means that the observation Z does not give us enough information to reliably distinguish between θ and θ' , but the consequences of mistaking one for the other are severe (since $d(\theta, \theta')$ is large).

Another useful consequence of Lemma 1 comes about when we consider the Bayesian setting, i.e., when we have a *prior distribution* π on Θ and consider the *average risk* of any estimator $\hat{\theta}$:

$$\mathbb{E}_{\pi} [d(\theta, \hat{\theta}(Z))] \triangleq \int_{\Theta} \pi(d\theta) \mathbb{E}_{\theta} [d(\theta, \hat{\theta}(Z))].$$

Then we have the following lower bound on $\mathfrak{M}(\Theta)$:

$$\mathfrak{M}(\Theta) \geq \inf_{\hat{\theta}} \mathbb{E}_{\pi} [d(\theta, \hat{\theta}(Z))].$$

To prove this, we simply note that, for any $\hat{\theta}$,

$$\mathbb{E}_{\pi} [d(\theta, \hat{\theta}(Z))] \triangleq \int_{\Theta} \pi(d\theta) \mathbb{E}_{\theta} [d(\theta, \hat{\theta}(Z))] \leq \sup_{\theta \in \Theta} \mathbb{E}_{\theta} [d(\theta, \hat{\theta})].$$

(In fact, under suitable regularity conditions, one can prove that the minimax risk is equal to the *worst-case Bayes risk*:

$$\mathfrak{M}(\Theta) = \inf_{\hat{\theta}} \sup_{\pi \in \mathcal{P}(\Theta)} \mathbb{E}_{\pi} [d(\theta, \hat{\theta}(Z))] = \sup_{\pi \in \mathcal{P}(\Theta)} \inf_{\hat{\theta}} \mathbb{E}_{\pi} [d(\theta, \hat{\theta}(Z))],$$

where the first equality always holds, while the second equality is valid whenever the conditions of the minimax theorem are satisfied.) With these definitions, we have the following:²

Corollary 2. *Let π be any prior distribution on Θ , and let μ be any joint probability distribution of a random pair $(\theta, \theta') \in \Theta \times \Theta$, such that the marginal distributions of both θ and θ' are equal to π . Then*

$$\mathfrak{M}(\Theta) \geq \inf_{\hat{\theta}} \mathbb{E}_{\pi} [d(\theta, \hat{\theta}(Z))] \geq \mathbb{E}_{\mu} [d(\theta, \theta') \cdot (1 - \|P_{\theta} - P_{\theta'}\|_{\text{TV}})]. \quad (16)$$

²I have learned this particular formulation from my colleagues Yihong Wu (here at Illinois) and Yury Polyanskiy (MIT).

Proof. Take the expectation of both sides of (13) with respect to μ and then use the fact that, under μ , both θ and θ' have the same distribution π . \square

In many cases, the analysis of minimax lower bounds can be reduced to the following problem: Suppose that the parameter set Θ can be identified with the m -dimensional binary hypercube $\{0, 1\}^m$ for some m , and the objective is to determine every bit of the underlying unknown parameter $\theta \in \{0, 1\}^m$ based on an observation $Z \sim P_\theta$. In that setting, a key result known as *Assouad's lemma* says that the difficulty of estimating the *entire* bit string θ is related to the difficulty of estimating each bit of θ separately, assuming you already know all other bits:

Theorem 3 (Assouad's lemma). *Suppose that $\Theta = \{0, 1\}^m$ and consider the Hamming metric*

$$d_H(\theta, \theta') \triangleq \sum_{i=1}^m 1_{\{\theta_i \neq \theta'_i\}}. \quad (17)$$

Then

$$\mathfrak{M}(\Theta) \geq \frac{m}{2} \left(1 - \max_{\theta, \theta': d_H(\theta, \theta')=1} \|P_\theta - P_{\theta'}\|_{\text{TV}} \right). \quad (18)$$

Proof. Notice that we can write $d_H(\theta, \theta')$ as a sum $\sum_{i=1}^m d_i(\theta, \theta')$, where $d_i(\theta, \theta') \triangleq 1_{\{\theta_i \neq \theta'_i\}}$, and each d_i is a pseudometric. Now let π be the uniform distribution on $\Theta = \{0, 1\}^m$, and for each $i \in \{1, \dots, m\}$ let μ_i be the joint distribution of a random pair $(\theta, \theta') \in \Theta \times \Theta$, under which $\theta \sim \pi$ and θ' differs from θ only in the i th bit, i.e., $\theta \sim \pi$ and $d_i(\theta, \theta') = 1$. Then the marginal distribution of θ' under μ_i is

$$\sum_{\theta \in \{0, 1\}^m} \mu_i(\theta, \theta') = \frac{1}{2^m} \sum_{\theta \in \{0, 1\}^m} 1_{\{\theta_i \neq \theta'_i \text{ and } \theta_j = \theta'_j, j \neq i\}} = \frac{1}{2^m} = \pi(\theta'),$$

since for each θ' there is only one θ that differs from it in a single bit. Now the idea is to apply Corollary 2 separately to each bit:

$$\begin{aligned} \mathfrak{M}(\Theta) &= \inf_{\hat{\theta}} \mathbb{E}_\pi [d_H(\theta, \hat{\theta}(Z))] \\ &= \inf_{\hat{\theta}} \sum_{i=1}^m \mathbb{E}_\pi [d_i(\theta, \hat{\theta}(Z))] \\ &\geq \sum_{i=1}^m \inf_{\hat{\theta}} \mathbb{E}_\pi [d_i(\theta, \hat{\theta}(Z))] \\ &\geq \sum_{i=1}^m \frac{1}{2} \mathbb{E}_{\mu_i} [d_i(\theta, \theta') \cdot (1 - \|P_\theta - P_{\theta'}\|_{\text{TV}})], \end{aligned}$$

where the last step is by Corollary 2. Since $d_i(\theta, \theta') = 1$ under μ_i for every i , we have

$$\begin{aligned} \mathfrak{M}(\Theta) &\geq \frac{1}{2} \sum_{i=1}^m \mathbb{E}_{\mu_i} [1 - \|P_\theta - P_{\theta'}\|_{\text{TV}}] \\ &\geq \frac{1}{2} \sum_{i=1}^m \min_{\theta, \theta': d_H(\theta, \theta')=1} (1 - \|P_\theta - P_{\theta'}\|_{\text{TV}}) \\ &= \frac{m}{2} \left(1 - \max_{\theta, \theta': d_H(\theta, \theta')=1} \|P_\theta - P_{\theta'}\|_{\text{TV}} \right), \end{aligned}$$

and we are done. \square

In order to apply Assouad's lemma, we need to have an upper bound on the total variation distance between any two P_θ and $P_{\theta'}$ with $d_H(\theta, \theta') = 1$. More often than not, it is easier to work with another distance between probability distributions, the so-called *Hellinger distance*. For any two $P, Q \in \mathcal{P}(Z)$, their *squared* Hellinger distance is given by

$$H^2(P, Q) \triangleq \frac{1}{2} \sum_{z \in Z} \left(\sqrt{P(z)} - \sqrt{Q(z)} \right)^2.$$

The total variation distance can be both upper- and lower-bounded by Hellinger:

$$\frac{1}{2} H^2(P, Q) \leq \|P - Q\|_{\text{TV}} \leq H(P, Q). \quad (19)$$

Moreover, for any n pairs of distributions $(P_1, Q_1), \dots, (P_n, Q_n)$,

$$H^2(P_1 \times \dots \times P_n, Q_1 \times \dots \times Q_n) \leq \sum_{i=1}^n H^2(P_i, Q_i). \quad (20)$$

Armed with these facts, we can prove the following version of Assouad's lemma:

Corollary 3. *Let $\{P_\theta : \theta \in \{0, 1\}^m\}$ be a collection of probability distributions on some set Z indexed by the vertices of the binary hypercube $\Theta = \{0, 1\}^m$. Suppose that there exists some constant $\alpha > 0$, such that*

$$H^2(P_\theta, P_{\theta'}) \leq \alpha, \quad \text{if } d_H(\theta, \theta') = 1. \quad (21)$$

Consider the problem of estimating the parameter $\theta \in \Theta$ based on n i.i.d. observations from P_θ , where the loss is measured by the Hamming distance. Then the corresponding minimax risk, which we denote by $\mathfrak{M}_n(\Theta)$, is lower-bounded as

$$\mathfrak{M}_n(\Theta) \geq \frac{m}{2} (1 - \sqrt{\alpha n}). \quad (22)$$

Proof. For any $\theta \in \Theta$, let P_θ^n denote the product of n copies of P_θ , i.e., the joint distribution of n i.i.d. samples from P_θ . For any two $\theta, \theta' \in \Theta$ with $d_H(\theta, \theta') = 1$, we have

$$\begin{aligned} \|P_\theta^n - P_{\theta'}^n\|_{\text{TV}} &\leq H(P_\theta^n, P_{\theta'}^n) \\ &\leq \sqrt{\sum_{i=1}^n H^2(P_\theta, P_{\theta'})} \\ &\leq \sqrt{\alpha n}, \end{aligned}$$

where the first step is by (19), the second is by (20), and the third is by (21). The bound (22) then follows from Assouad's lemma. \square

Assouad's lemma is a wonderful tool, but it is only applicable when the parameter space Θ is a binary hypercube and the metric d is the Hamming metric. To handle other situations, we need different methods. Let's recall what we had seen earlier: the minimax risk $\mathfrak{M}(\Theta)$ is large if we can find at least two distinct points $\theta, \theta' \in \Theta$, such that (a) the distance $d(\theta, \theta')$ is large and (b) the distributions P_θ and $P_{\theta'}$ are statistically close to one another, i.e., the total variation distance $\|P_\theta - P_{\theta'}\|_{\text{TV}}$ is small. Let's see what happens when we follow this logic further and consider any two $\theta, \theta' \in \Theta$ that are separated from each other by a distance of at least $\varepsilon > 0$:

Lemma 2. Fix $\varepsilon > 0$, and let $\Theta_0 \subset \Theta$ be any set which is ε -separated, i.e., $d(\theta, \theta') \geq \varepsilon$ for any two distinct $\theta, \theta' \in \Theta_0$. Then

$$\mathfrak{M}(\Theta) \geq \frac{\varepsilon}{2} \inf_{\tilde{\theta}} \sup_{\theta \in \Theta_0} \mathbb{P}_{\theta}(\tilde{\theta} \neq \theta), \quad (23)$$

where the infimum is over all estimators $\tilde{\theta}$ that take values in Θ_0 .

Proof. Fix any ε -separated set Θ_0 . Then, since $\Theta_0 \subset \Theta$, we have $\mathfrak{M}(\Theta) \geq \mathfrak{M}(\Theta_0)$. Given any estimator $\hat{\theta} = \hat{\theta}(Z)$ taking values in Θ , define another estimator $\tilde{\theta} = \tilde{\theta}(Z)$ taking values in Θ_0 as follows:

$$\tilde{\theta} = \operatorname{argmin}_{\theta' \in \Theta_0} d(\hat{\theta}, \theta').$$

Then, for any $\theta \in \Theta_0$,

$$d(\tilde{\theta}, \theta) \leq d(\tilde{\theta}, \hat{\theta}) + d(\hat{\theta}, \theta) \leq 2d(\hat{\theta}, \theta),$$

where the first step is by the triangle inequality, while the second step is by definition of $\tilde{\theta}$. Consequently,

$$\mathfrak{M}(\Theta) \geq \mathfrak{M}(\Theta_0) \geq \frac{1}{2} \inf_{\tilde{\theta}} \sup_{\theta \in \Theta_0} \mathbb{E}_{\theta} [d(\tilde{\theta}, \theta)]. \quad (24)$$

Now let us consider the event $\{z \in Z : \tilde{\theta}(z) \neq \theta\}$ for some $\theta \in \Theta_0$. Since $\tilde{\theta} \in \Theta_0$ by construction, and Θ_0 is ε -separated, we have $\{z \in Z : \tilde{\theta}(z) \neq \theta\} \subseteq \{z \in Z : d(\tilde{\theta}(z), \theta) \geq \varepsilon\}$. Using this fact and Markov's inequality, we can write

$$\mathbb{P}_{\theta}(\tilde{\theta} \neq \theta) \leq \mathbb{P}_{\theta}(d(\tilde{\theta}, \theta) \geq \varepsilon) \leq \frac{\mathbb{E}_{\theta}[d(\tilde{\theta}, \theta)]}{\varepsilon}. \quad (25)$$

Using this in (24), we get (23). \square

Lemma 2 is at its most effective when we can find a *finite* set Θ_0 which is ε -separated and has large cardinality. The reason for the latter requirement is that then it should not be too hard to arrange things in such a way that the probability of error $\mathbb{P}_{\theta}(\tilde{\theta} \neq \theta)$ is uniformly bounded away from zero for *any* estimator $\tilde{\theta}$. In particular, if we can find an ε -separated set $\Theta_0 \subset \Theta$, such that

$$\inf_{\tilde{\theta}} \max_{\theta \in \Theta_0} \mathbb{P}_{\theta}(\tilde{\theta}, \theta) \geq \alpha$$

for some absolute constant $\alpha > 0$, then we can lower-bound the minimax risk $\mathfrak{M}(\Theta)$ by $\alpha\varepsilon/2$. So now we need tight lower bounds on the probability of error when testing multiple hypotheses that hold for *any* estimation procedure. Bounds of this sort rely, in one way or another, on a fundamental result from information theory known as *Fano's inequality* [CT06] (or Fano's lemma, as statisticians call it [Yu97]). We will use a particular version, due to Lucien Birgé [Bir05]. To state it, we first need to define relative entropy (or Kullback–Leibler divergence). Given any two probability distributions P, Q on Z , the *relative entropy* (or *Kullback–Leibler divergence*) between P and Q is defined as

$$D(P\|Q) \triangleq \begin{cases} \sum_{z \in Z} P(z) \log \frac{P(z)}{Q(z)}, & \text{if } \operatorname{supp}(P) \subseteq \operatorname{supp}(Q) \\ +\infty, & \text{otherwise} \end{cases} \quad (26)$$

where $\operatorname{supp}(P) \triangleq \{z \in Z : P(z) > 0\}$ is the support of P . Here are some useful properties:

- $D(P\|Q) \geq 0$, with equality if and only if $P \equiv Q$;
- for any n pairs of distributions (P_i, Q_i) , $1 \leq i \leq n$, we have

$$D(P_1 \times \dots \times P_n \| Q_1 \times \dots \times Q_n) = \sum_{i=1}^n D(P_i \| Q_i). \quad (27)$$

Also, one can bound the total variation distance in terms of the KL divergence as follows:

$$\|P - Q\|_{\text{TV}} \leq \sqrt{\frac{1}{2} D(P\|Q)}. \quad (28)$$

This inequality is usually referred to as *Pinsker's inequality* after Mark S. Pinsker, although Pinsker didn't prove it. In the present form, the inequality (28), with the optimal constant 1/2 in front of $D(P\|Q)$, was obtained independently by I. Csiszár, J.H.B. Kemperman, and S. Kullback.

With these preliminaries out of the way, we can state Birgé's bound:

Lemma 3 (Birgé). *Let $\{P_i\}_{i=0}^N$ be a finite family of probability distributions on Z , and let $\{A_i\}_{i=0}^N$ be a family of disjoint events. Then*

$$\min_{0 \leq i \leq N} P_i(A_i) \leq \max\left(\alpha, \frac{\bar{K}}{\log(N+1)}\right),$$

where $\alpha = 0.71$ and

$$\bar{K} = \frac{1}{N} \sum_{i=1}^N D(P_i \| P_0).$$

Now let's see how we can apply Birgé's lemma: Let $\Theta_0 = \{\theta_0, \dots, \theta_N\}$ be some ε -separated subset of Θ . Fix an estimator $\tilde{\theta}$ taking values in Θ_0 . For each $0 \leq i \leq N$, consider the probability distribution $P_i = P_{\theta_i}$ and the event $A_i = \{z \in Z : \tilde{\theta}(z) = \theta_i\}$. Since the elements of Θ_0 are all distinct, the events A_0, \dots, A_N are disjoint. Now suppose that we have chosen Θ_0 in such a way that

$$\bar{K} = \frac{1}{N} \sum_{i=1}^N D(P_0 \| P_i) = \frac{1}{N} \sum_{i=1}^N D(P_{\theta_0} \| P_{\theta_i}) \leq \alpha \log(N+1). \quad (29)$$

Then, by Birgé's lemma,

$$\begin{aligned} \max_{\theta \in \Theta_0} \mathbb{P}_\theta(\tilde{\theta} \neq \theta) &= \max_{0 \leq i \leq N} P_i(A_i^c) \\ &= 1 - \min_{0 \leq i \leq N} P_i(A_i) \\ &\geq 1 - \alpha \\ &\equiv 0.29. \end{aligned}$$

This bound holds for *any* estimator $\tilde{\theta}$. Consequently, using Lemma 2, we get the lower bound $\mathfrak{M}(\Theta) \geq (1 - \alpha)\varepsilon/2$.

2 Proof of Theorem 1

Roughly speaking, the strategy of the proof will be as follows: We start by observing that we can lower-bound the minimax risk $R_n(h, \mathcal{F})$ by

$$R_n(h, \mathcal{F}) \geq \inf_{\tilde{f}_n} \sup_{P \in \mathcal{Q}} \mathbb{E} [L(\tilde{f}_n) - L(f^*)]$$

for any subset $\mathcal{Q} \subset \mathcal{P}(h, \mathcal{F})$. We will then construct \mathcal{Q} in such a way that its elements can be naturally indexed by the vertices of a binary hypercube of dimension $V-1$ (where V is the VC dimension of \mathcal{F}), and the expected excess risk for each $P \in \mathcal{Q}$ can be related to the minimax risk for the problem of estimating the corresponding binary string of length $V-1$. At that point, we will be in a position to apply Assouad's lemma.

2.1 Construction of \mathcal{Q}

As stated above, our set \mathcal{Q} of “difficult” distributions will consist of 2^{V-1} distributions P_b , $b \in \{0, 1\}^{V-1}$. To construct these distributions, we will first pick the marginal distribution $P_X \in \mathcal{P}(X)$ of the feature X , and then specify the conditional distributions $P_{Y|X}^{(b)}$ of the binary label Y given X , for each $b \in \{0, 1\}^{V-1}$. For now, let us assume that $h > 0$.

We construct P_X as follows. Since \mathcal{F} is a VC class with VC dimension V , there exists a set $\{x_1, \dots, x_V\} \subset X$ of points that can be shattered by \mathcal{F} , i.e., for any binary string $\beta \in \{0, 1\}^V$ there exists at least one $f \in \mathcal{F}$, such that $f(x_j) = \beta_j$ for all $1 \leq j \leq V$. Given a parameter $p \in [0, 1/(V-1)]$, which will be chosen later, we choose P_X so that $P_X(\{x_1, \dots, x_V\}) = 1$, and

$$P_X(x_j) = \begin{cases} p, & \text{if } 1 \leq j \leq V-1 \\ 1 - (V-1)p, & \text{otherwise} \end{cases}. \quad (30)$$

In other words, P_X is a discrete distribution supported on the shattered set $\{x_1, \dots, x_V\}$ that puts the same mass p on each of the first $V-1$ points and dumps the rest of the probability mass, equal to $1 - (V-1)p$, on the last point x_V . Owing to our condition on p , this is a valid probability distribution.

Next we choose $P_{Y|X}^{(b)}$ for each $b \in \{0, 1\}^{V-1}$ as follows: For a fixed b , the conditional distribution of Y given $X = x$ is

$$\begin{aligned} & \text{Bernoulli}\left(\frac{1 + (2b_j - 1)h}{2}\right), & \text{if } x = x_j \text{ and } j \in \{1, \dots, V-1\} \\ & \text{Bernoulli}(0), & \text{otherwise} \end{aligned}$$

In other words, for a fixed b , we have

$$\eta_b(x) \equiv P_{Y|X}^{(b)}(Y = 1 | X = x) = \begin{cases} \frac{1-h}{2}, & \text{if } x = x_j \text{ for some } j \in \{1, \dots, V-1\}, \text{ and } b_j = 0 \\ \frac{1+h}{2}, & \text{if } x = x_j \text{ for some } j \in \{1, \dots, V-1\}, \text{ and } b_j = 1. \\ 0, & \text{otherwise} \end{cases} \quad (31)$$

Therefore, the corresponding Bayes classifier, which we denote by f_b^* , is given by

$$f_b^*(x) = \begin{cases} 0, & \text{if } x = x_j \text{ for some } j \in \{1, \dots, V-1\}, \text{ and } b_j = 0 \\ 1, & \text{if } x = x_j \text{ for some } j \in \{1, \dots, V-1\}, \text{ and } b_j = 1. \\ 0, & \text{otherwise} \end{cases} \quad (32)$$

That is, the output of the Bayes classifier on each x_j , $1 \leq j \leq V-1$, is simply equal to the bit value b_j , and is zero everywhere else.

It remains to check whether the resulting set $\mathcal{Q} = \{P_b : b \in \{0, 1\}^{V-1}\}$ is contained in $\mathcal{P}(h, \mathcal{F})$. First of all, from (31) we see that $|2\eta_b(x) - 1| \geq h$ for all x (indeed, $|2\eta_b(x) - 1| = h$ when $x \in \{x_1, \dots, x_{V-1}\}$, and $|2\eta_b(x) - 1| = 1$ otherwise). Secondly, because $\{x_1, \dots, x_V\}$ is shattered by \mathcal{F} , there exists at least one $f \in \mathcal{F}$, such that $f_b^*(x) = f(x)$ for all $x \in \{x_1, \dots, x_V\}$. Thus, $\mathcal{Q} \subset \mathcal{P}(h, \mathcal{F})$, as claimed.

2.2 Reduction to an estimation problem on the binary hypercube

Now that we have constructed \mathcal{Q} , we will show that the problem of learning a classifier when the n training samples are drawn i.i.d. from some $P_b \in \mathcal{Q}$ is at least as difficult as reconstructing the underlying bit string b — indeed, intuitively this makes sense, given the structure of the Bayes classifier f_b^* corresponding to P_b .

We start by noting that, for any classifier $f : X \rightarrow \{0, 1\}$ and any distribution P on $X \times \{0, 1\}$, we have³

$$L(f) - L(f^*) = \mathbb{E} [|2\eta(X) - 1| |f(X) - f^*(X)|].$$

From this, we see that if $P \in \mathcal{P}(h, \mathcal{F})$, then

$$L(f) - L(f^*) \geq h \mathbb{E} [|f(X) - f^*(X)|] \equiv h \|f - f^*\|_{L_1},$$

where the L_1 norm is computed w.r.t. the marginal distribution P_X . Hence, we have

$$\begin{aligned} \inf_{\tilde{f}_n} \sup_{P \in \mathcal{Q}} \mathbb{E} [L(\tilde{f}_n) - L(f^*)] &= \inf_{\tilde{f}_n} \max_{b \in \{0, 1\}^{V-1}} \mathbb{E}_b [L(\tilde{f}_n) - L(f_b^*)] \\ &\geq h \inf_{\tilde{f}_n} \max_{b \in \{0, 1\}^{V-1}} \mathbb{E}_b \|\tilde{f}_n - f_b^*\|_{L_1}, \end{aligned} \quad (33)$$

where the L_1 norm is w.r.t. the distribution P_X defined in (30), and $\mathbb{E}_b[\cdot]$ denotes expectation w.r.t. P_b . Now we are ready to carry out the reduction to an estimation problem on the binary hypercube $\{0, 1\}^{V-1}$. Given any candidate estimator \tilde{f}_n , let \tilde{b}_n be an estimator that takes values in $\{0, 1\}^{V-1}$, and is defined as follows:

$$\tilde{b}_n \triangleq \operatorname{argmin}_{b \in \{0, 1\}^{V-1}} \|\tilde{f}_n - f_b^*\|_{L_1}. \quad (34)$$

In other words, \tilde{b}_n is the binary string that indexes the element of $\{f_b^* : b \in \{0, 1\}^{V-1}\}$ which is the closest to \tilde{f}_n in the L_1 norm. Then, for any b ,

$$\begin{aligned} \|f_{\tilde{b}_n}^* - f_b^*\|_{L_1} &\leq \|f_{\tilde{b}_n}^* - \tilde{f}_n\|_{L_1} + \|\tilde{f}_n - f_b^*\|_{L_1} \\ &\leq 2\|\tilde{f}_n - f_b^*\|_{L_1}, \end{aligned} \quad (35)$$

where the first step is by the triangle inequality, while the second step is by (34). Combining (35) with (33), we get

$$\inf_{\tilde{f}_n} \sup_{P \in \mathcal{Q}} \mathbb{E} [L(\tilde{f}_n) - L(f^*)] \geq \frac{h}{2} \inf_{\tilde{b}_n} \max_{b \in \{0, 1\}^{V-1}} \mathbb{E}_b \|f_{\tilde{b}_n}^* - f_b^*\|_{L_1}, \quad (36)$$

³Exercise: prove this!

where the infimum is over all estimators that take values in $\{0, 1\}^{V-1}$ based on n i.i.d. samples from one of the P_b 's.

Now let us inspect the L_1 norm $\|f_b^* - f_{b'}^*\|_{L_1}$ for any two b, b' . Using (32), we have

$$\begin{aligned}\|f_b^* - f_{b'}^*\|_{L_1} &= \int_{\mathcal{X}} |f_b^*(x) - f_{b'}^*(x)| P_X(dx) \\ &= \sum_{j=1}^V P_X(x_j) |f_b^*(x_j) - f_{b'}^*(x_j)| \\ &= p \sum_{j=1}^{V-1} |f_b^*(x_j) - f_{b'}^*(x_j)| \\ &= p \sum_{j=1}^{V-1} |b_j - b'_j| \\ &= p d_H(b, b').\end{aligned}$$

Substituting this into (36) gives

$$\inf_{\tilde{f}_n} \sup_{P \in \mathcal{Q}} \mathbb{E} [L(\tilde{f}_n) - L(f^*)] \geq \frac{ph}{2} \inf_{\hat{b}_n} \max_{b \in \{0, 1\}^{V-1}} \mathbb{E}_b [d_H(\hat{b}_n, b)], \quad (37)$$

as claimed. Now we are in a position to apply Assouad's lemma.

2.3 Applying Assouad's lemma

In order to apply Assouad's lemma, we need an upper bound on the squared Hellinger distance $H^2(P_b, P_{b'})$ for all b, b' with $d_H(b, b') = 1$. This is a simple computation. In fact, for any two $b, b' \in \{0, 1\}^{V-1}$ we have

$$\begin{aligned}H^2(P_b, P_{b'}) &= \sum_{j=1}^V \sum_{y \in \{0, 1\}} \left(\sqrt{P_b(x_j, y)} - \sqrt{P_{b'}(x_j, y)} \right)^2 \\ &= p \sum_{j=1}^{V-1} \sum_{y \in \{0, 1\}} \left(\sqrt{P_{Y|X}^{(b)}(y|x_j)} - \sqrt{P_{Y|X}^{(b')}(y|x_j)} \right)^2 \\ &= p \sum_{j=1}^{V-1} H^2 \left(\text{Bernoulli} \left(\frac{1 + (2b_j - 1)h}{2} \right), \text{Bernoulli} \left(\frac{1 + (2b'_j - 1)h}{2} \right) \right).\end{aligned}$$

For each $j \in \{1, \dots, V-1\}$, the j th term in the above summation is nonzero if and only if $b_j \neq b'_j$, in which case it is equal to the squared Hellinger distance between the $\text{Bernoulli}(\frac{1-h}{2})$ and $\text{Bernoulli}(\frac{1+h}{2})$ distributions. Thus,

$$\begin{aligned}H^2(P_b, P_{b'}) &= p H^2 \left(\text{Bernoulli} \left(\frac{1-h}{2} \right), \text{Bernoulli} \left(\frac{1+h}{2} \right) \right) \cdot d_H(b, b') \\ &= 2p \left(\sqrt{\frac{1-h}{2}} - \sqrt{\frac{1+h}{2}} \right)^2 d_H(b, b') \\ &= 2p (1 - \sqrt{1-h^2}) d_H(b, b').\end{aligned}$$

In particular, the collection of distributions $\{P_b : b \in \{0, 1\}^{V-1}\}$ satisfies the condition (21) of Corollary 3 with $\alpha = 2p(1 - \sqrt{1 - h^2}) \leq 2ph^2$. Therefore, applying the bound (22), we get

$$\inf_{\tilde{f}_n} \sup_{P \in \mathcal{Q}} \mathbb{E} [L(\tilde{f}_n) - L(f^*)] \geq \frac{p(V-1)h}{4} \left(1 - \sqrt{2nph^2}\right).$$

If we let $p = 2/9nh^2$, the term in parentheses will be equal to $1/3$, and

$$\inf_{\tilde{f}_n} \sup_{P \in \mathcal{Q}} \mathbb{E} [L(\tilde{f}_n) - L(f^*)] \geq \frac{V-1}{54nh},$$

assuming that the condition $p \leq 1/(V-1)$ holds. This will be the case if $h \geq \sqrt{(V-1)/n}$. Therefore,

$$\inf_{\tilde{f}_n} \sup_{P \in \mathcal{P}(h, \mathcal{F})} \mathbb{E} [L(\tilde{f}_n) - L(f^*)] \geq \frac{V-1}{54nh}, \quad \text{if } h \geq \sqrt{(V-1)/n}. \quad (38)$$

If $h \leq \sqrt{(V-1)/n}$, we can use the above construction with $\tilde{h} = \sqrt{(V-1)/n}$. In that case, because $\mathcal{P}(h, \mathcal{F}) \subseteq \mathcal{P}(\tilde{h}, \mathcal{F})$ whenever $h \geq \tilde{h}$, we see that

$$\begin{aligned} \inf_{\tilde{f}_n} \sup_{P \in \mathcal{P}(h, \mathcal{F})} \mathbb{E} [L(\tilde{f}_n) - L(f^*)] &\geq \inf_{\tilde{f}_n} \sup_{P \in \mathcal{P}(\tilde{h}, \mathcal{F})} \mathbb{E} [L(\tilde{f}_n) - L(f^*)] \\ &\geq \frac{V-1}{54n\tilde{h}} \\ &= \frac{1}{54} \sqrt{\frac{V-1}{n}}, \quad \text{if } h \leq \sqrt{(V-1)/n}. \end{aligned} \quad (39)$$

Putting together (38) and (39), we get the bound of Theorem 1.

3 Proof of Theorem 2

In its broad outline, the proof is very similar to the proof of Theorem 1. Fix some $N \geq 4D$. Since \mathcal{F} is (N, D) -rich, there exist N distinct points $x_1, \dots, x_N \in \mathcal{X}$, such that

$$\{0, 1\}_D^N \triangleq \{b \in \{0, 1\}^N : \text{wt}(b) = D\} \subseteq \mathcal{F}(x^N),$$

where $\text{wt}(b)$ denotes the Hamming weight of a binary string b . Let P_X be the uniform distribution on $\{x_1, \dots, x_N\}$. Also, for each $b \in \{0, 1\}_D^N$, let

$$\eta_b(x_i) \triangleq \frac{1 + (2b_i - 1)h}{2}, \quad 1 \leq i \leq N.$$

Now for each such b define a distribution $P_b \in \mathcal{P}(\mathcal{X} \times \{0, 1\})$ by $P_b = P_X \times P_{Y|X}^{(b)}$, where

$$P_{Y|X}^{(b)}(1|x_i) = \frac{1 + (2b_i - 1)h}{2} \equiv \eta_b(x_i).$$

In other words, the conditional distribution $P_{Y|X=x_i}^{(b)}$ is a Bernoulli measure with bias $\frac{1+h}{2}$ if $b_i = 1$ and $\frac{1-h}{2}$ if $b_i = 0$. The corresponding Bayes classifier is given by $f_b^*(x_i) = b_i$ for each i , as before. Moreover,

since \mathcal{F} is (N, D) -rich, for any $b \in \{0, 1\}_D^N$ there exists at least one $f \in \mathcal{F}$, such that $f(x_i) = b_i = f_b^*(x_i)$ for every i . Consequently,

$$\mathcal{Q} = \{P_b : b \in \{0, 1\}_D^N\} \subset \mathcal{P}(h, \mathcal{F}).$$

Proceeding in the same way as in the proof of Theorem 1, for any subset $\mathcal{C} \subseteq \{0, 1\}_D^N$,

$$\begin{aligned} R_n(h, \mathcal{F}) &\geq \frac{h}{2} \inf_{\tilde{b}_n \in \mathcal{C}} \max_{b \in \mathcal{C}} \mathbb{E}_b \left\| f_{\tilde{b}_n}^* - f_b^* \right\|_{L_1} \\ &= \frac{h}{2N} \inf_{\tilde{b}_n \in \mathcal{C}} \max_{b \in \mathcal{C}} \mathbb{E}_b [d_{\text{H}}(\tilde{b}_n, b)]. \end{aligned}$$

Since $\{0, 1\}_D^N$ is not the entire binary hypercube $\{0, 1\}^N$, we cannot use Assouad's lemma. Instead, we will use Lemma 2 in conjunction with Birgé's bound. In order to do that, we must first construct a nice, well-separated subset of $\{0, 1\}_D^N$. The following combinatorial result, due to P. Reynaud-Bouret, does the trick. There exists a subset $\mathcal{C} \subset \{0, 1\}_D^N$ with the following properties:

1. $d_{\text{H}}(b, b') > D/2$ for any two distinct $b, b' \in \mathcal{C}$;
2. $\log|\mathcal{C}| \geq \kappa D \log \frac{N}{D}$, where $\kappa \approx 0.233$.

Item 1 says that this \mathcal{C} is $D/2$ -separated in the Hamming distance. Thus, by Lemma 2,

$$\begin{aligned} R_n(h, \mathcal{F}) &\geq \frac{hD}{4N} \inf_{\tilde{b}_n \in \mathcal{C}} \max_{b \in \mathcal{C}} \mathbb{P}_b(\tilde{b}_n \neq b) \\ &= \frac{hD}{4N} \left(1 - \sup_{\tilde{b}_n \in \mathcal{C}} \min_{b \in \mathcal{C}} \mathbb{P}_b(\tilde{b}_n \neq b) \right). \end{aligned}$$

We are now in a position to apply Birgé's lemma:

$$\sup_{\tilde{b}_n \in \mathcal{C}} \min_{b \in \mathcal{C}} \mathbb{P}_b(\tilde{b}_n \neq b) \leq \alpha,$$

where $\alpha = 0.71$, provided

$$\bar{K} = \frac{1}{|\mathcal{C}| - 1} \sum_{b \in \mathcal{C}, b \neq b_0} D(P_b^n \| P_{b_0}^n) \leq \alpha \log|\mathcal{C}|, \quad (40)$$

where b_0 is some fixed but arbitrary element of \mathcal{C} , and P_b^n is a product of n copies of P_b (recall that our estimators are based on an i.i.d. sample of size n). So now we turn to the analysis of \bar{K} . For any two $b, b' \in \{0, 1\}^N$, we have

$$\begin{aligned} D(P_b^n \| P_{b'}^n) &= nD(P_b \| P_{b'}) \\ &= n \sum_{i=1}^N \sum_{y \in \{0, 1\}} P_b(x_i, y) \log \frac{P_b(x_i, y)}{P_{b'}(x_i, y)} \\ &= \frac{n}{N} \sum_{i=1}^N \sum_{y \in \{0, 1\}} P_{Y|X}^{(b)}(y|x_i) \log \frac{P_{Y|X}^{(b)}(y|x_i)}{P_{Y|X}^{(b')}(y|x_i)} \\ &= \frac{n}{N} \sum_{i=1}^N D\left(\text{Bernoulli}\left(\frac{1 + (2b_i - 1)h}{2}\right) \parallel \text{Bernoulli}\left(\frac{1 + (2b'_i - 1)h}{2}\right)\right), \end{aligned}$$

where the first line is by (27), while the rest follows from definitions. Now, in the above summation, the i th term is nonzero if and only if $b_i = b'_i$, and in that case it is equal to

$$D\left(\text{Bernoulli}\left(\frac{1+h}{2}\right)\left\|\text{Bernoulli}\left(\frac{1-h}{2}\right)\right.\right) = h \log \frac{1+h}{1-h}.$$

Consequently,

$$\begin{aligned} D(P_b^n \| P_{b'}^n) &= \frac{n}{N} h \log \left(\frac{1+h}{1-h}\right) \sum_{i=1}^N 1_{\{b_i \neq b'_i\}} \\ &= \frac{n}{N} h \log \left(\frac{1+h}{1-h}\right) \cdot d_H(b, b'). \end{aligned}$$

Now, for any two $b, b' \in \{0, 1\}_D^N$, $d_H(b, b') \leq \text{wt}(b) + \text{wt}(b') = 2D$. Moreover, using the bound $\log t \leq t - 1$ for any $t \geq 0$, we have

$$h \log \frac{1+h}{1-h} \leq h \left(\frac{1+h}{1-h} - 1\right) = \frac{2h^2}{1-h}.$$

Therefore, for any two $b, b' \in \mathcal{C}$,

$$D(P_b^n \| P_{b'}^n) \leq \frac{4nh^2D}{N(1-h)},$$

which implies that $\bar{K} \leq \frac{4nh^2D}{N(1-h)}$. Using this and the fact that $\log |\mathcal{C}| \leq \kappa D \log \frac{N}{D}$ with $\kappa \approx 0.233$, we see that the condition (40) will be satisfied if we can choose some $N \geq 4D$, so that

$$N \log \frac{N}{D} \geq \frac{4nh^2}{\alpha\kappa(1-h)}. \quad (41)$$

A tedious calculation (see [MN06]) shows that the choice

$$N = \left\lceil \frac{8nh^2}{\alpha\kappa(1-h) \left(1 + \log \frac{nh^2}{D}\right)} \right\rceil$$

does the job, provided $h \geq \sqrt{D/n}$. This completes the proof.

References

- [Bir05] L. Birgé. A new lower bound for multiple hypothesis testing. *IEEE Transactions on Information Theory*, 51(4):1611–1615, April 2005.
- [CT06] T. M. Cover and J. T. Thomas. *Elements of Information Theory*. Wiley, 2nd edition, 2006.
- [MN06] P. Massart and É. Nédélec. Risk bounds for statistical learning. *Annals of Statistics*, 34(5):2326–2366, 2006.
- [Vap98] V. Vapnik. *Statistical Learning Theory*. Wiley, 1998.
- [Yu97] B. Yu. Assouad, Fano, and Le Cam. In D. Pollard, E. Torgersen, and G. Yang, editors, *Festschrift for Lucien Le Cam*, pages 423–435. Springer, 1997.