

Fat-Shattering and the Learnability of Real-Valued Functions

PETER L. BARTLETT

Department of Systems Engineering, Research School of Information Sciences and Engineering, Australian National University, Canberra, 0200 Australia

PHILIP M. LONG

Department of Information Systems and Computer Sciences, National University of Singapore, Singapore 119260, Republic of Singapore

AND

ROBERT C. WILLIAMSON

Department of Engineering, Australian National University, Canberra, 0200 Australia

Received February 3, 1994

We consider the problem of learning real-valued functions from random examples when the function values are corrupted with noise. With mild conditions on independent observation noise, we provide characterizations of the learnability of a real-valued function class in terms of a generalization of the Vapnik–Chervonenkis dimension, the fat-shattering function, introduced by Kearns and Schapire. We show that, given some restrictions on the noise, a function class is learnable in our model if and only if its fat-shattering function is finite. With different (also quite mild) restrictions, satisfied for example by gaussian noise, we show that a function class is learnable from polynomially many examples if and only if its fat-shattering function grows polynomially. We prove analogous results in an agnostic setting, where there is no assumption of an underlying function class. © 1996 Academic Press, Inc.

1. INTRODUCTION

In many common definitions of learning, a learner sees a sequence of values of an unknown function at random points, and must, with high probability, choose an accurate approximation to that function. The function is assumed to be a member of some known class. Using a popular definition of the problem of learning $\{0, 1\}$ -valued functions (probably approximately correct learning—see [12, 26]), Blumer, Ehrenfeucht, Haussler, and Warmuth have shown [12] that the Vapnik–Chervonenkis dimension (see [27]) of a function class characterizes its learnability, in the sense that a function class is learnable if and only if its Vapnik–Chervonenkis dimension is finite. Natarajan [19] and Ben-David, Cesa-Bianchi, Haussler, and Long [11] have characterized the learnability of $\{0, \dots, n\}$ -valued functions for fixed n . Alon, Ben-David, Cesa-Bianchi, and Haussler have proved an analogous result for the problem of learning probabilistic concepts [1]. In this case, there is

an unknown $[0, 1]$ -valued function, but the learner does not receive a sequence of values of the function at random points. Instead, with each random point it sees either 0 or 1, with the probability of a 1 given by the value of the unknown function at that point. Kearns and Schapire [16] introduced a generalization of the Vapnik–Chervonenkis dimension, which we call the fat-shattering function, and showed that a class of probabilistic concepts is learnable only if the class has a finite fat-shattering function. The main learning result of [1] is that finiteness of the fat-shattering function of a class of probabilistic concepts is also sufficient for learnability.

In this paper, we consider the learnability of $[0, 1]$ -valued function classes. We show that a class of $[0, 1]$ -valued functions is learnable from a finite training sample with observation noise satisfying some mild conditions (the distribution has bounded support and its density satisfies a smoothness constraint) if and only if the class has a finite fat-shattering function. Here, as elsewhere, our main contribution is in showing that the finiteness of the fat-shattering function is necessary for learning. We also consider small-sample learnability, for which the sample size is allowed to grow only polynomially with the required performance parameters. We show that a real-valued function class is learnable from a small sample with observation noise satisfying some other quite mild conditions (the distribution need not have bounded support, but it must have light tails and be symmetric about zero; gaussian noise satisfies these conditions) if and only if the fat-shattering function of the class has a polynomial rate of growth. We also consider agnostic learning [15, 17], in which there is no assumption of an underlying function generating the training examples, and the performance of the learning algorithm is measured by comparison with some function

class F . We show that the fat-shattering function of F characterizes finite-sample and small-sample learnability in this case also. In fact, the proof in [1] that finiteness of the fat-shattering function of a class of probabilistic concepts implies learnability also gives a related sufficient condition for agnostic learnability of $[0, 1]$ -valued functions. We show that this condition is implied by finiteness of the fat-shattering function of F .

The proof of the lower bound on the number of examples necessary for learning is in two steps. First, we show that the problem of learning real-valued functions in the presence of noise is not much easier than that of learning functions in a discrete-valued function class obtained by quantizing the real-valued function class. This formalizes the intuition that a noisy, real-valued measurement provides little more information than a quantized measurement, if the quantization width is sufficiently small. Existing lower bounds on the number of examples required for learning discrete-valued function classes [11, 19] are not strong enough for our purposes. We improve these lower bounds by relating the problem of learning the quantized function class to that of learning $\{0, 1\}$ -valued functions.

In addition to the aforementioned papers, other general results about learning real-valued functions have been obtained. Haussler [15] gives sufficient conditions for agnostic learnability. Anthony, Bartlett, Ishai, and Shawe-Taylor [4] give necessary and sufficient conditions that a function that approximately interpolates the target function is a good approximation to it (see also [5, 3]). Natarajan [20] considers the problem of learning a class of real-valued functions in the presence of bounded observation noise and presents sufficient conditions for learnability. (Theorem 2 in [4] shows that these conditions are not necessary in our setting.) Merhav and Feder [18], and Auer, Long, Maass, and Woeginger [6] study function learning in a worst-case setting.

In the next section, we define admissible noise distribution classes and the learning problems and present the characterizations of learnability. Sections 3 and 4 give lower and upper bounds on the number of examples necessary for learning real-valued functions. Section 5 presents the characterization of agnostic learnability. Section 6 discusses our results. An earlier version of this paper appeared in [10].

2. DEFINITIONS AND MAIN RESULT

Denote the integers by \mathbb{Z} , the positive integers by \mathbb{N} , the reals by \mathbb{R} and the nonnegative reals by \mathbb{R}^+ . We use \log to denote logarithm to base 2, and \ln to denote the natural logarithm. Fix an arbitrary set X . Throughout the paper, X denotes the input space on which the real-valued functions are defined. We refer to probability distributions on X without explicitly defining a σ -algebra \mathcal{S} . For countable X ,

let \mathcal{S} be the set of all subsets of X . If X is a metric space, let \mathcal{S} be the Borel sets of X . All functions and sets we consider are assumed to be measurable.

Classes of Noise Distributions

The noise distributions we consider are absolutely continuous, and their densities have bounded variation. A function $f: \mathbb{R} \rightarrow \mathbb{R}$ is said to have **bounded variation** if there is a constant $C > 0$ such that for every ordered sequence $x_0 < \dots < x_n$ in \mathbb{R} ($n \in \mathbb{N}$) we have

$$\sum_{k=1}^n |f(x_k) - f(x_{k-1})| \leq C.$$

In that case, the **total variation** of f on \mathbb{R} is

$$V(f) = \sup \left\{ \sum_{k=1}^n |f(x_k) - f(x_{k-1})| : n \in \mathbb{N}, x_0 < \dots < x_n \right\}.$$

DEFINITION 1. An **admissible noise distribution class** \mathcal{D} is a class of distributions on \mathbb{R} that satisfies

1. Each distribution in \mathcal{D} has mean 0 and finite variance.
2. Each distribution in \mathcal{D} is absolutely continuous and its probability density function (pdf) has bounded variation. Furthermore, there is a function $v: \mathbb{R}^+ \rightarrow \mathbb{R}^+$ such that, if f is the pdf of any distribution in \mathcal{D} with variance σ^2 , then $V(f) \leq v(\sigma)$. The function v is called the **total variation function** of the class \mathcal{D} .

If \mathcal{D} also satisfies the following condition, we say it is a **bounded admissible noise distribution class**.

3. There is a function $s: \mathbb{R}^+ \rightarrow \mathbb{R}^+$ such that, if D is a distribution in \mathcal{D} with variance σ^2 , then the support of D is contained in a closed interval of length $s(\sigma)$. The function s is called the **support function** of \mathcal{D} .

If \mathcal{D} satisfies Conditions 1, 2, and the following condition,¹ we say it is an **almost-bounded admissible noise distribution class**.

- 3'. Each distribution D in \mathcal{D} has an even pdf ($f(x) = f(-x)$) and light tails: there are constants s_0 and c_0 in \mathbb{R}^+ such that, for all distributions D in \mathcal{D} with variance σ^2 and all $s > s_0\sigma$,

$$D\{\eta: |\eta| > s/2\} \leq c_0 e^{-s/\sigma}.$$

¹ In fact, Condition 3' is stronger than we need. It suffices that the distributions be "close to" symmetric and have light tails in the following sense: there are constants s_0 and c_0 in \mathbb{R}^+ such that, for all distributions D in \mathcal{D} with variance σ^2 , and all $s > s_0\sigma$, if $l \in \mathbb{R}$ satisfies $\int_{l+s}^{l+s} x\phi(x) dx = 0$, then $\int_{l+s}^{l+s} \phi(x) dx \geq 1 - c_0 e^{-s/\sigma}$, where ϕ is the pdf of D .

EXAMPLE (Uniform noise). Let $\mathcal{U} = \{U_\sigma : \sigma > 0\}$, where U_σ is uniform on $(-\sqrt{3}\sigma, \sqrt{3}\sigma)$. Then this noise has mean 0, standard deviation σ , total variation function $v(\sigma) = 1/(\sqrt{3}\sigma)$, and support function $s(\sigma) = 2\sqrt{3}\sigma$, so \mathcal{U} is a bounded admissible noise distribution class.

EXAMPLE (Gaussian noise). Let $\mathcal{G} = \{G_\sigma : \sigma > 0\}$, where G_σ is the zero mean gaussian distribution with variance σ^2 . Since the density f_σ of G_σ has $f_\sigma(0) = (\sqrt{2\pi}\sigma)^{-1}$, and $f_\sigma(x)$ is monotonic decreasing for $x > 0$, the total variation function is $v(\sigma) = 2(\sqrt{2\pi}\sigma)^{-1}$. Obviously, f_σ is an even function. Standard bounds on the area under the tails of the gaussian density (see [21, p. 64, Fact 3.7.3]) give

$$G_\sigma\{\eta \in \mathbb{R} : |\eta| > s/2\} \leq \exp\left(-\frac{s^2}{8\sigma^2}\right), \quad (1)$$

and if $s > 8\sigma$, $\exp(-s^2/(8\sigma^2)) < \exp(-s/\sigma)$, so the constants $c_0 = 1$ and $s_0 = 8$ will satisfy Condition 3'. So the class \mathcal{G} of gaussian distributions is almost-bounded admissible.

The Learning Problem

Choose a set F of functions from X to $[0, 1]$. For $m \in \mathbb{N}$, $f \in F$, $x \in X^m$, and $\eta \in \mathbb{R}^m$, let

$$\begin{aligned} \text{sam}(x, \eta, f) \\ = ((x_1, f(x_1) + \eta_1), \dots, (x_m, f(x_m) + \eta_m)) \in (X \times \mathbb{R})^m. \end{aligned}$$

(We often dispense with the parentheses in tuples of this form, to avoid cluttering the notation.) Informally, a **learning algorithm** takes a sample of the above form and outputs a hypothesis for f . More formally, a **deterministic learning algorithm**² is defined to be a mapping from $\bigcup_m (X \times \mathbb{R})^m$ to $[0, 1]^X$. A **randomized learning algorithm** L is a pair (A, P_Z) , where P_Z is a distribution on a set Z , and A is a mapping from $\bigcup_m (X \times \mathbb{R})^m \times Z^m$ to $[0, 1]^X$. That is, given a sample of length m , the randomized algorithm chooses a sequence $z \in Z^m$ at random from P_Z^m and passes it to the (deterministic) mapping A as a parameter.

For a probability distribution P on X , $f \in F$, and $h : X \rightarrow [0, 1]$, define

$$\mathbf{er}_{P,f}(h) = \int_X |h(x) - f(x)| dP(x).$$

The following definition of learning is based on those of [12, 19, 26].

DEFINITION 2. Let \mathcal{D} be a class of distributions on \mathbb{R} . Choose $0 < \varepsilon$, $\delta < 1$, $\sigma > 0$, and $m \in \mathbb{N}$. We say a learning

² Despite the name ‘‘algorithm,’’ there is no requirement that this mapping be computable. Throughout the paper, we ignore issues of computability.

algorithm $L = (A, P_Z)$ $(\varepsilon, \delta, \sigma)$ -**learns** F from m examples **with noise** \mathcal{D} if for all distributions P on X , all functions f in F , and all distributions $D \in \mathcal{D}$ with variance σ^2 ,

$$\begin{aligned} (P^m \times D^m \times P_Z^m)\{(x, \eta, z) \in X^m \times \mathbb{R}^m \times Z^m : \\ \mathbf{er}_{P,f}(A(\text{sam}(x, \eta, f), z)) \geq \varepsilon\} < \delta. \end{aligned}$$

Similarly, L (ε, δ) -**learns** F from m examples **without noise** if, for all distributions P on X and all functions f in F ,

$$\begin{aligned} (P^m \times P_Z^m)\{(x, z) \in X^m \times Z^m : \mathbf{er}_{P,f}(A(\text{sam}(x, 0, f), z)) \geq \varepsilon\} \\ < \delta. \end{aligned}$$

We say F is **learnable with noise** \mathcal{D} if there is a learning algorithm L and a function $m_0 : (0, 1) \times (0, 1) \times \mathbb{R}^+ \rightarrow \mathbb{N}$ such that for all $0 < \varepsilon$, $\delta < 1$, for all $\sigma > 0$, algorithm $L(\varepsilon, \delta, \sigma)$ learns F from $m_0(\varepsilon, \delta, \sigma)$ examples with noise \mathcal{D} . We say F is **small-sample learnable with noise** \mathcal{D} if, in addition, the function m_0 is bounded by a polynomial in $1/\varepsilon$, $1/\delta$, and σ .

The following definition comes from [16]. Choose $x_1, \dots, x_d \in X$. We say x_1, \dots, x_d are γ -**shattered by** F if there exists $r \in [0, 1]^d$ such that for each $b \in \{0, 1\}^d$, there is an $f \in F$ such that for each i

$$f(x_i) \begin{cases} \geq r_i + \gamma & \text{if } b_i = 1 \\ \leq r_i - \gamma & \text{if } b_i = 0. \end{cases}$$

For each γ , let

$$\text{fat}_F(\gamma) = \max\{d \in \mathbb{N} : \exists x_1, \dots, x_d, F \text{ } \gamma\text{-shatters } x_1, \dots, x_d\}$$

if such a maximum exists, and ∞ otherwise. If $\text{fat}_F(\gamma)$ is finite for all γ , we say F has a **finite fat-shattering function**.

The following is our main result.

THEOREM 3. *Suppose F is a permissible³ class of $[0, 1]$ -valued functions defined on X . If \mathcal{D} is a bounded admissible noise distribution class, then F is learnable with observation noise \mathcal{D} if and only if F has a finite fat-shattering function. If \mathcal{D} is an almost-bounded admissible noise distribution class, then F is small-sample learnable with observation noise \mathcal{D} if and only if there is a polynomial p such that $\text{fat}_F(\gamma) < p(1/\gamma)$ for all $\gamma > 0$.*

3. LOWER BOUND

In this section, we give a lower bound on the number of examples necessary to learn a real-valued function class in the presence of observation noise. Lemma 5 in Section 3.1

³ This is a benign measurability constraint defined in Section 4.

shows that an algorithm that can learn a real-valued function class with observation noise can be used to construct an algorithm that can learn a quantized version of the function class to slightly worse accuracy and confidence with the same number of examples, provided the quantization width is sufficiently small. Lemma 10 in Section 3.2 gives a lower bound on the number of examples necessary for learning a quantized function class in terms of its fat-shattering function. In Section 3.3, we combine these results to give the lower bound for real-valued functions, Theorem 11.

3.1. Learnability with Noise Implies Quantized Learnability

In this subsection, we relate the problem of learning a real-valued function class with observation noise to the problem of learning a quantized version of that class, without noise.

DEFINITION 4. For $\alpha \in \mathbb{R}^+$, define the quantization function

$$Q_\alpha(y) = \alpha \left\lceil \frac{y - \alpha/2}{\alpha} \right\rceil.$$

For a set $S \subset \mathbb{R}$, let $Q_\alpha(S) = \{Q_\alpha(y) : y \in S\}$. For a function class $F \subset [0, 1]^X$, let $Q_\alpha(F)$ be the set $\{Q_\alpha \circ f : f \in F\}$ of $Q_\alpha([0, 1])$ -valued functions defined on X .

LEMMA 5. Suppose F is a set of functions from X to $[0, 1]$, \mathcal{D} is an admissible noise distribution class with total variation function v , A is a learning algorithm, $0 < \varepsilon, \delta < 1$, $\sigma \in \mathbb{R}^+$, and $m \in \mathbb{N}$. If the quantization width $\alpha \in \mathbb{R}^+$ satisfies

$$\alpha \leq \min\left(\frac{\delta}{v(\sigma)m}, 2\varepsilon\right),$$

and $A(\varepsilon, \delta, \sigma)$ -learns F from m examples with noise \mathcal{D} , then there is a randomized learning algorithm (C, P_Z) that $(2\varepsilon, 2\delta)$ -learns $Q_\alpha(F)$ from m examples.

Figure 1 illustrates our approach. Suppose an algorithm A can $(\varepsilon, \delta, \sigma)$ -learn from m noisy examples $(x_i, f(x_i) + \eta_i)$. If we quantize the observations to accuracy α and add noise that is uniform on $(-\alpha/2, \alpha/2)$, Lemma 6(a) shows that the distribution of the observations is approximately unchanged (in the notation of Fig. 1, the distributions P_1 and P_2 are close), so A learns almost as well as it did previously. If we define Algorithm B as this operation of adding uniform noise and then invoking Algorithm A , B solves a quantized learning problem in which the examples are given as $(x_i, Q_\alpha(f(x_i) + \eta_i))$. Lemma 6(b) shows that this problem is similar to the problem of learning the quantized function class when the observations are contaminated with independent noise whose distribution is a quantized version of the original observation noise (that is, the examples

are given as $(x_i, Q_\alpha(f(x_i)) + Q_\alpha(v_i))$). In the notation of Fig. 1, Lemma 6(b) shows that the distributions P_3 and P_4 are close. It follows that Algorithm C , which adds this quantized noise to the observations and passes them to Algorithm B , learns the quantized function class without observation noise (that is, when the examples are given as $(x_i, Q_\alpha(f(x_i)))$).

For distributions P and Q on \mathbb{R} , define the **total variation distance** between P and Q as

$$d_{\text{TV}}(P, Q) = 2 \sup_E |P(E) - Q(E)|,$$

where the supremum is over all Borel sets. If P and Q are discrete, it is easy to show that

$$d_{\text{TV}}(P, Q) = \sum_x |P(x) - Q(x)|,$$

where the sum is over all x in the union of the supports of P and Q . Similarly, if P and Q are continuous with probability density functions p and q , respectively,

$$d_{\text{TV}}(P, Q) = \int_{-\infty}^{\infty} |p(x) - q(x)| dx.$$

LEMMA 6. Let \mathcal{D} be an admissible noise distribution class with total variation function v . Let $\sigma > 0$ and $0 < \alpha < 1$. Let D be a distribution in \mathcal{D} with variance σ^2 . Let η, ζ , and v be random variables, and suppose that η and v are distributed according to D and ζ is distributed uniformly on $(-\alpha/2, \alpha/2)$.

(a) For any $y \in [0, 1]$, if P_1 is the distribution of $y + \eta$ and P_2 is the distribution of $Q_\alpha(y + \eta) + \zeta$, we have

$$d_{\text{TV}}(P_1, P_2) \leq \alpha v(\sigma).$$

(b) For any $y \in [0, 1]$, if P_3 is the distribution of $Q_\alpha(y + \eta)$ and P_4 is the distribution of $Q_\alpha(y) + Q_\alpha(v)$, we have

$$d_{\text{TV}}(P_3, P_4) \leq \alpha v(\sigma).$$

Proof. Let p be the pdf of D .

(a) The random variable $y + \eta$ has density $p_1(a) = p(a - y)$, and $Q_\alpha(y + \eta) + \zeta$ has density p_2 given by

$$p_2(a) = \frac{1}{\alpha} \int_{Q_\alpha(a) - \alpha/2}^{Q_\alpha(a) + \alpha/2} p(x - y) dx$$

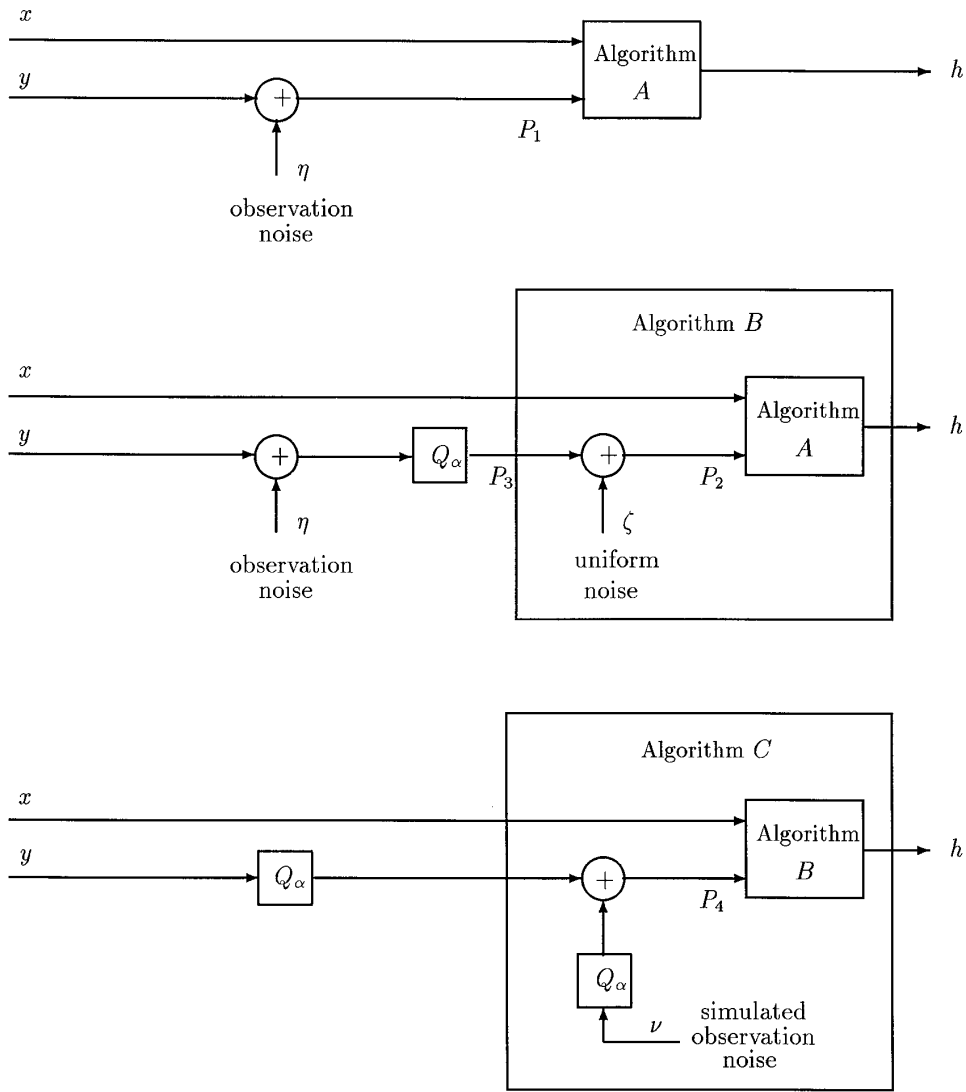


FIG. 1. Lemma 5 shows that a learning algorithm for real-valued functions (Algorithm A) can be used to construct a randomized learning algorithm for quantized functions (Algorithm C).

for $a \in \mathbb{R}$. So

$$\begin{aligned}
 & d_{\text{TV}}(P_1, P_2) \\
 &= \int_{-\infty}^{\infty} \left| p(x-y) - \frac{1}{\alpha} \int_{Q_\alpha(x)-\alpha/2}^{Q_\alpha(x)+\alpha/2} p(\theta-y) d\theta \right| dx \\
 &= \sum_{n=-\infty}^{\infty} \int_{-\alpha/2}^{\alpha/2} \left| p(x-y+n\alpha) - \frac{1}{\alpha} \int_{-\alpha/2}^{\alpha/2} p(\theta-y+n\alpha) d\theta \right| dx \\
 &= \int_{-\alpha/2}^{\alpha/2} \sum_{n=-\infty}^{\infty} \left| p(x-y+n\alpha) - \frac{1}{\alpha} \int_{-\alpha/2}^{\alpha/2} p(\theta-y+n\alpha) d\theta \right| dx.
 \end{aligned}$$

By the mean value theorem, there are z_1 and z_2 in $[-\alpha/2, \alpha/2]$ such that

$$p(z_1 - y + n\alpha) \leq \frac{1}{\alpha} \int_{-\alpha/2}^{\alpha/2} p(\theta - y + n\alpha) d\theta \leq p(z_2 - y + n\alpha),$$

so for all $x \in [-\alpha/2, \alpha/2]$,

$$\begin{aligned}
 & \sum_{n=-\infty}^{\infty} \left| p(x-y+n\alpha) - \frac{1}{\alpha} \int_{-\alpha/2}^{\alpha/2} p(\theta-y+n\alpha) d\theta \right| \\
 & \leq \sum_{n=-\infty}^{\infty} \sup_{z \in (-\alpha/2, \alpha/2)} |p(x-y+n\alpha) - p(z-y+n\alpha)| \\
 & \leq v(\sigma),
 \end{aligned}$$

and therefore,

$$d_{\text{TV}}(P_1, P_2) \leq \alpha v(\sigma).$$

(b) The distribution P_3 of $Q_\alpha(y + \eta)$ is discrete and is given by

$$P_3(a) = \begin{cases} \int_{n\alpha - \alpha/2}^{n\alpha + \alpha/2} p(x - y) dx & \text{if } a = n\alpha \text{ for some } n \in \mathbb{Z} \\ 0 & \text{otherwise.} \end{cases}$$

Since v has distribution D , the distribution P_4 of $Q_\alpha(y) + Q_\alpha(v)$ is also discrete and is given by

$$P_4(a) = \begin{cases} \int_{n\alpha - \alpha/2}^{n\alpha + \alpha/2} p(x) dx & \text{if } a = Q_\alpha(y) + n\alpha \text{ for some } n \in \mathbb{Z} \\ 0 & \text{otherwise.} \end{cases}$$

So

$$\begin{aligned} d_{\text{TV}}(P_3, P_4) &= \sum_{n=-\infty}^{\infty} |P_3(n\alpha) - P_4(n\alpha)| \\ &= \sum_{n=-\infty}^{\infty} \left| \int_{n\alpha - \alpha/2}^{n\alpha + \alpha/2} p(x - y) dx - \int_{n\alpha - \alpha/2}^{n\alpha + \alpha/2} p(x - Q_\alpha(y)) dx \right| \\ &\leq \sum_{n=-\infty}^{\infty} \int_{-\alpha/2}^{\alpha/2} |p(x - y + n\alpha) - p(x - Q_\alpha(y) + n\alpha)| dx \\ &= \int_{-\alpha/2}^{\alpha/2} \sum_{n=-\infty}^{\infty} |p(x - y + n\alpha) - p(x - Q_\alpha(y) + n\alpha)| dx \\ &\leq \int_{-\alpha/2}^{\alpha/2} \sum_{n=-\infty}^{\infty} \sup_{z \in (-\alpha/2, \alpha/2)} |p(x - y + n\alpha) - p(x - y + n\alpha + z)| dx \\ &\leq \alpha v(\sigma). \quad \blacksquare \end{aligned}$$

We will use the following lemma. The proof is by induction and is implicit in the proof of Lemma 12 in [8].

LEMMA 7. *If P_i and Q_i ($i = 1, \dots, m$) are distributions on a set Y , and χ is a $[0, 1]$ -valued random variable defined on Y^m , then*

$$\left| \int_{Y^m} \chi dP - \int_{Y^m} \chi dQ \right| \leq \frac{1}{2} \sum_{i=1}^m d_{\text{TV}}(P_i, Q_i),$$

where $P = \prod_{i=1}^m P_i$ and $Q = \prod_{i=1}^m Q_i$ are distributions on Y^m .

Proof (of Lemma 5). We will describe a randomized algorithm (Algorithm C) that is constructed from Algo-

rithm A and show that it $(2\varepsilon, 2\delta)$ -learns the quantized function class $Q_\alpha(F)$. Fix a noise distribution D in \mathcal{D} with variance σ^2 , a function $f \in F$, and a distribution P on X . Since A $(\varepsilon, \delta, \sigma)$ -learns F , we have

$$P^m \times D^m \{ (x, \eta) \in X^m \times \mathbb{R}^m : \mathbf{er}_{P,f}(A(\text{sam}(x, \eta, f))) \geq \varepsilon \} < \delta.$$

That is, the probability (over all $x \in X^m$ and $\eta \in \mathbb{R}^m$) that Algorithm A chooses a bad function is small. We will show that this implies that the probability that Algorithm C chooses a bad function is also small, where the probability is over all $x \in X^m$ and all values of the random variables that Algorithm C uses.

Let ζ be a random variable with distribution U_α , where U_α is the uniform distribution on $(-\alpha/2, \alpha/2)$. For an arbitrary sequence (y_1, \dots, y_m) , let Algorithm B be the randomized algorithm that adds noise ζ to each y value it receives and passes the sequence to Algorithm A . That is, for any sequence of (x_i, y_i) pairs,

$$B(x_1, y_1, \dots, x_m, y_m) = A(x_1, y_1 + \zeta_1, \dots, x_m, y_m + \zeta_m).$$

First we prove that, for a given sequence x of input values, the probability that Algorithm A outputs a bad hypothesis when it is called from Algorithm B in the scenario shown in Fig. 1 (that is, when it sees examples of the form $(x_i, Q_\alpha(f(x_i) + \eta_i) + \zeta_i)$) is no more than $\delta/2$ more than the probability that Algorithm A outputs a bad hypothesis after receiving examples $(x_i, f(x_i) + \eta_i)$. We prove this by considering the set of noisy function values for the input sequence x that cause Algorithm A to output a bad hypothesis.

Now, fix a sequence $x = (x_1, \dots, x_m) \in X^m$, and define the events

$$E = \{ \eta \in \mathbb{R}^m : \mathbf{er}_{P,f}(A(\text{sam}(x, \eta, f))) \geq \varepsilon \},$$

$$E_1 = \{ y \in \mathbb{R}^m : \mathbf{er}_{P,f}(A(x_1, y_1, \dots, x_m, y_m)) \geq \varepsilon \}.$$

That is, E is the set of noise sequences that make A choose a bad function, and E_1 is the corresponding set of y sequences. Clearly,

$$D^m(E) = \left(\prod_{i=1}^m P_{1|x_i} \right) (E_1), \quad (2)$$

where $P_{1|x_i}$ is the distribution of $f(x_i) + \eta$. We will show that $D^m(E)$ is close to the corresponding probability under the distribution of y values that Algorithm A sees when Algorithm B invokes it.

Define $P_{2|x_i}$ as the distribution of $Q_\alpha(f(x_i) + \eta) + \zeta$. From Lemma 6a, $d_{\text{TV}}(P_{1|x_i}, P_{2|x_i}) \leq \alpha v(\sigma)$. Applying Lemma 7 with $\chi = 1_{E_1}$, the indicator⁴ function for E_1 , gives

$$\left(\prod_{i=1}^m P_{2|x_i} \right) (E_1) - \left(\prod_{i=1}^m P_{1|x_i} \right) (E_1) \leq m\alpha v(\sigma)/2.$$

By hypothesis $\alpha \leq \delta/(mv(\sigma))$, so this and (2) imply

$$\left(\prod_{i=1}^m P_{2|x_i} \right) (E_1) \leq D^m(E) + \delta/2.$$

Next we observe that, for a fixed sequence x of input values, the probability that Algorithm B outputs a bad hypothesis when given quantized noisy examples (of the form $(x_i, Q_\alpha(f(x_i) + \eta_i))$) is equal to the probability that Algorithm A outputs a bad hypothesis when given examples of the form $(x_i, Q_\alpha(f(x_i) + \eta_i) + \zeta_i)$. More formally, we can write this as follows. Let $P_{3|x_i}$ be the distribution of $Q_\alpha(f(x_i) + \eta)$, and let

$$E_3 = \{(y, \zeta) \in \mathbb{R}^m \times \mathbb{R}^m: \\ \mathbf{er}_{P, f}(A(x_1, y_1 + \zeta_1, \dots, x_m, y_m + \zeta_m)) \geq \varepsilon\}.$$

In this case, E_3 is the set of (y, ζ) pairs that correspond to B choosing a bad function. Clearly,

$$\left(\prod_{i=1}^m P_{2|x_i} \right) (E_1) = \left(\prod_{i=1}^m P_{3|x_i} \times U_\alpha^m \right) (E_3).$$

Let v be a random variable with distribution D . Let Algorithm C be the randomized algorithm that adds noise $Q_\alpha(v)$ to each y value it receives, and passes the sequence to Algorithm B . That is,

$$C(x_1, y_1, \dots, x_m, y_m) \\ = B(x_1, y_1 + Q_\alpha(v_1), \dots, x_m, y_m + Q_\alpha(v_m)).$$

Next we prove that, for a fixed sequence x , the probability that Algorithm B outputs a bad hypothesis when it is called from Algorithm C in the scenario shown in Fig. 1 (that is, when it sees examples of the form $(x_i, Q_\alpha(f(x_i) + Q_\alpha(v_i)))$) is no more than $\delta/2$ more than the probability that Algorithm B outputs a bad hypothesis after receiving examples of the form $(x_i, Q_\alpha(f(x_i) + \eta_i))$.

Let $P_{4|x_i}$ be the distribution of $Q_\alpha(f(x_i)) + Q_\alpha(v)$. Applying Lemma 7, with χ equal the probability under U_α that A produces a bad hypothesis, gives

$$\left| \left(\prod_{i=1}^m P_{3|x_i} \times U_\alpha^m \right) (E_3) - \left(\prod_{i=1}^m P_{4|x_i} \times U_\alpha^m \right) (E_3) \right| \\ \leq \sum_{i=1}^m d_{\text{TV}}(P_{3|x_i}, P_{4|x_i})/2.$$

From Lemma 6(b), $d_{\text{TV}}(P_{4|x_i}, P_{3|x_i}) \leq \alpha v(\sigma)$, so we have

$$\left(\prod_{i=1}^m P_{4|x_i} \times U_\alpha^m \right) (E_3) \leq \left(\prod_{i=1}^m P_{3|x_i} \times U_\alpha^m \right) (E_3) + \delta/2 \\ = \left(\prod_{i=1}^m P_{2|x_i} \right) (E_1) + \delta/2 \\ \leq D^m(E) + \delta.$$

We have shown that, for any $x \in X^m$,

$$U_\alpha^m \times D^m\{(\zeta, v): \mathbf{er}_{P, f}(A(x_1, Q_\alpha(f(x_1)) + Q_\alpha(v_1) + \zeta_1, \dots, \\ x_m, Q_\alpha(f(x_m)) + Q_\alpha(v_m) + \zeta_m)) \geq \varepsilon\} \\ \leq D^m\{\eta: \mathbf{er}_{P, f}(A(\text{sam}(x, \eta, f))) \geq \varepsilon\} + \delta.$$

It follows that

$$P^m \times U_\alpha^m \times D^m\{(x, \zeta, v): \mathbf{er}_{P, f}(A(x_1, Q_\alpha(f(x_1)) \\ + Q_\alpha(v_1) + \zeta_1, \dots, x_m, Q_\alpha(f(x_m)) + Q_\alpha(v_m) + \zeta_m)) \geq \varepsilon\} \\ \leq P^m \times D^m\{(x, \eta): \mathbf{er}_{P, f}(A(\text{sam}(x, \eta, f))) \geq \varepsilon\} + \delta \\ \leq 2\delta.$$

For any function $h: X \rightarrow [0, 1]$, the triangle inequality for the absolute difference on \mathbb{R} gives

$$\mathbf{er}_{P, f}(h) \\ = \int_X |h(x) - f(x)| dP(x) \\ \geq \int_X (|h(x) - Q_\alpha(f(x))| - |f(x) - Q_\alpha(f(x))|) dP(x) \\ \geq \mathbf{er}_{P, Q_\alpha(f)}(h) - \alpha/2 \\ \geq \mathbf{er}_{P, Q_\alpha(f)}(h) - \varepsilon,$$

since for all x , $|f(x) - Q_\alpha(f(x))| \leq \alpha/2$, and $\alpha \leq 2\varepsilon$ by hypothesis. It follows that

$$\{\mathbf{er}_{P, Q_\alpha(f)}(C(\dots)) \geq 2\varepsilon\} \subseteq \{\mathbf{er}_{P, f}(C(\dots)) \geq \varepsilon\}.$$

⁴ That is, $1_{E_1}(y)$ takes value 1 if $y \in E_1$, and 0 otherwise.

Hence

$$\Pr(\mathbf{er}_{P, Q_\alpha(f)}(C(x_1, Q_\alpha(f(x_1)), \dots, x_m, Q_\alpha(f(x_m)))) \geq 2\varepsilon) < 2\delta,$$

where the probability is taken over all x in X^m and all values of ζ and v , the random variables used by Algorithm C . This is true for any $Q_\alpha(f)$ in $Q_\alpha(F)$, so this algorithm $(2\varepsilon, 2\delta)$ -learns $Q_\alpha(F)$ from m examples. ■

3.2. Lower Bounds for Quantized Learning

In the previous subsection, we showed that if a class F can be $(\varepsilon, \delta, \sigma)$ -learned with a certain number of examples; then an associated class $Q_\alpha(F)$ of discrete-valued functions can be $(2\varepsilon, 2\delta)$ -learned with the same number of examples. Given this result, one would be tempted to apply techniques of Natarajan [19] or Ben-David, Cesa-Bianchi, Haussler, and Long [11] (who consider the learnability of discrete-valued functions) to lower bound the number of examples required for learning $Q_\alpha(F)$. The main results of those papers, however, were for the discrete loss function, where the learner “loses” 1 whenever its hypothesis is incorrect. When those results are applied directly to get bounds for learning with the absolute loss, the resulting bounds are not strong enough for our purposes because of the restrictions on α required to show that learning F is not much harder than learning $Q_\alpha(F)$.

In this subsection, we present a new technique, inspired by the techniques of [7]. We show that an algorithm for learning a class of discrete-valued functions can effectively be used as a subroutine in an algorithm for learning binary-valued functions. We then apply a lower bound result for binary-valued functions.

For each $d \in \mathbb{N}$, let POWER_d be the set of all functions from $\{1, \dots, d\}$ to $\{0, 1\}$. We will make use of the following special case of a general result about POWER_d [12, Theorem 2.1(b)].

THEOREM 8 [12]. *Let A be a randomized learning algorithm which always outputs $\{0, 1\}$ -valued hypotheses. If A is given fewer than $d/2$ examples, A fails to $(1/8, 1/8)$ -learn POWER_d .*

Theorem 2.1(b) of [12] is stated for deterministic algorithms, but an almost identical proof gives the same result for randomized algorithms.

We will also make use of the standard Chernov bounds, proved in this form by Angluin and Valiant [2].

THEOREM 9 [2]. *Let Y_1, \dots, Y_m be independent, identically distributed $\{0, 1\}$ -valued random variables, where $\Pr(Y_1 = 1) = p$. Then*

$$\Pr\left(\sum_{i=1}^m Y_i \geq 2mp\right) \leq e^{-mp/3}$$

$$\Pr\left(\sum_{i=1}^m Y_i \leq mp/2\right) \leq e^{-mp/8}.$$

LEMMA 10. *For $0 < \alpha < \frac{1}{2}$, choose a set F of functions from X to $Q_\alpha([0, 1])$, $d \in \mathbb{N}$, and $\gamma > 0$ such that $\text{fat}_F(\gamma) \geq d$. If a randomized learning algorithm A is given fewer than*

$$\frac{d - 666}{4 + 192 \ln \lceil 1/\alpha + 1/2 \rceil}$$

examples, A fails to $(\gamma/32, 1/16)$ -learn F without noise.

Proof. We will show that if there is an algorithm that can $(\gamma/32, 1/16)$ -learn the quantized class F from fewer than the number of examples given in the lemma, then this could be used as a subroutine of an algorithm that could $(1/8, 1/8)$ -learn POWER_d from fewer than $d/2$ examples, violating Theorem 8.

Choose an algorithm A for learning F . Let $x_1, \dots, x_d \in X$ be γ -shattered by F , and let $r_1, \dots, r_d \in [0, 1]^d$ be such that for each $b \in \{0, 1\}^d$, there is an $f_b \in F$ such that for all j , $1 \leq j \leq d$,

$$f_b(x_j) \begin{cases} \geq r_j + \gamma & \text{if } b_j = 1 \\ \leq r_j - \gamma & \text{if } b_j = 0. \end{cases}$$

For each $q \in \mathbb{N}$, consider the algorithm \tilde{A}_q (which will be used for learning POWER_d) which uses A as a subroutine as follows. Given $m > q$ examples, $(\kappa_1, y_1), \dots, (\kappa_m, y_m)$ in $\{1, \dots, d\} \times \{0, 1\}$, Algorithm \tilde{A}_q first, for each $v \in Q_\alpha([0, 1])^q = \{0, \alpha, \dots, \alpha \lceil 1/\alpha - 1/2 \rceil\}^q$, sets $h_{\kappa, v} = A((x_{\kappa_1}, v_1), \dots, (x_{\kappa_q}, v_q))$. Algorithm \tilde{A}_q then uses this to define a set \tilde{S} of $\{0, 1\}$ -valued functions defined on $\{1, \dots, d\}$ by

$$\tilde{S} = \{\tilde{h}_{\kappa, v} : v \in Q_\alpha([0, 1])^q\},$$

where

$$\tilde{h}_{\kappa, v}(j) = \begin{cases} 1 & \text{if } h_{\kappa, v}(x_j) \geq r_j \\ 0 & \text{otherwise,} \end{cases}$$

for all $j \in \{1, \dots, d\}$. Finally, \tilde{A}_q returns an \tilde{h}^* in \tilde{S} for which the number of disagreements with the last $m - q$ examples is minimized. That is,

$$\tilde{h}^* = \arg \min_{\tilde{h} \in \tilde{S}} \{|\{j \in \{q+1, \dots, m\} : \tilde{h}(\kappa_j) \neq y_j\}|\}.$$

We claim that if A can $(\gamma/32, 1/16)$ -learn F from $m_0 \in \mathbb{N}$ examples without noise, then \tilde{A}_{m_0} can $(1/8, 1/8)$ -learn POWER_d from

$$m_0 + \lceil 96(\ln 32 + m_0 \ln \lceil 1/\alpha + 1/2 \rceil) \rceil$$

examples without noise, and we can then apply Theorem 8 to give the desired lower bound on m_0 . To see this, assume A $(\gamma/32, 1/16)$ -learns F from m_0 examples, and let $\tilde{A} = \tilde{A}_{m_0}$. Suppose \tilde{A} is trying to learn $g \in \text{POWER}_d$ and the distribution on the domain $\{1, \dots, d\}$ is \tilde{P} . Let P be the corresponding distribution on $\{x_1, \dots, x_d\}$, and let $b = (g(1), \dots, g(d)) \in \{0, 1\}^d$. Since A $(\gamma/32, 1/16)$ -learns F , we have

$$\begin{aligned} & \tilde{P}^{m_0}\{(\kappa_1, \dots, \kappa_{m_0}): \\ & \quad \mathbf{er}_{P, f_b}(A((x_{\kappa_1}, f_b(x_{\kappa_1})), \dots, (x_{\kappa_{m_0}}, f_b(x_{\kappa_{m_0}})))) \geq \gamma/32\} \\ & < 1/16, \end{aligned}$$

which implies

$$\begin{aligned} & \tilde{P}^{m_0}\{\kappa \in \{1, \dots, d\}^{m_0}: \forall v \in \mathcal{Q}_\alpha([0, 1])^{m_0}, \mathbf{er}_{P, f_b}(h_{\kappa, v}) \geq \gamma/32\} \\ & < 1/16. \end{aligned}$$

This can be rewritten as

$$\tilde{P}^{m_0}\left\{\kappa: \forall v, \int |h_{\kappa, v}(x_j) - f_b(x_j)| d\tilde{P}(j) \geq \gamma/32\right\} < 1/16,$$

which, applying Markov's inequality, yields

$$\tilde{P}^{m_0}\{\kappa: \forall v, \tilde{P}\{j: |h_{\kappa, v}(x_j) - f_b(x_j)| \geq \gamma\} \geq 1/32\} < 1/16. \quad (3)$$

Now, for all j , $|f_b(x_j) - r_j| \geq \gamma$; so if $|\tilde{h}_{\kappa, v}(j) - b_j| = 1$ the definitions of $\tilde{h}_{\kappa, v}$ and f_b imply $|h_{\kappa, v}(x_j) - f_b(x_j)| \geq \gamma$. Therefore, $\mathbf{er}_{\tilde{P}, g}(\tilde{h}_{\kappa, v}) \geq 1/32$ implies

$$\tilde{P}\{j: |h_{\kappa, v}(x_j) - f_b(x_j)| \geq \gamma\} \geq 1/32;$$

so (3) implies

$$\tilde{P}^{m_0}\{\kappa: \forall v, \mathbf{er}_{\tilde{P}, g}(\tilde{h}_{\kappa, v}) \geq 1/32\} < 1/16. \quad (4)$$

That is, \tilde{A} is unlikely to choose \tilde{S} so that all elements have large error. We will show that \tilde{A} can use the remaining u examples to find an accurate function in \tilde{S} . Let

$$u = \lceil 96(\ln 32 + m_0 \ln \lceil 1/\alpha + 1/2 \rceil) \rceil.$$

Fix a v in $\mathcal{Q}_\alpha([0, 1])^{m_0}$ and a κ in $\{1, \dots, d\}^{m_0}$. If $\mathbf{er}_{\tilde{P}, g}(\tilde{h}_{\kappa, v}) \geq 1/8$, we can apply Theorem 9, with

$$Y_j = \begin{cases} 1 & \text{if } \tilde{h}_{\kappa, v}(\lambda_j) \neq g(\lambda_j) \\ 0 & \text{otherwise,} \end{cases}$$

to give

$$\tilde{P}^u\{(\lambda_1, \dots, \lambda_u): |\{j: \tilde{h}_{\kappa, v}(\lambda_j) \neq g(\lambda_j)\}| \leq u/16\} \leq e^{-u/64}.$$

Similarly, if $\mathbf{er}_{\tilde{P}, g}(\tilde{h}_{\kappa, v}) \leq 1/32$, Theorem 9 implies

$$\tilde{P}^u\{(\lambda_1, \dots, \lambda_u): |\{j: \tilde{h}_{\kappa, v}(\lambda_j) \neq g(\lambda_j)\}| \geq u/16\} \leq e^{-u/96}.$$

Since this is true for any v and since $|\mathcal{Q}_\alpha([0, 1])| = \lceil 1/\alpha + 1/2 \rceil$, we have

$$\begin{aligned} & \tilde{P}^u\{(\lambda_1, \dots, \lambda_u): \exists v, (\mathbf{er}_{\tilde{P}, g}(\tilde{h}_{\kappa, v}) \geq 1/8 \\ & \quad \text{and } |\{j: \tilde{h}_{\kappa, v}(\lambda_j) \neq g(\lambda_j)\}| \leq u/16) \\ & \quad \text{or } (\mathbf{er}_{\tilde{P}, g}(\tilde{h}_{\kappa, v}) \leq 1/32 \\ & \quad \text{and } |\{j: \tilde{h}_{\kappa, v}(\lambda_j) \neq g(\lambda_j)\}| \geq u/16)\} \\ & \leq 2\lceil 1/\alpha + 1/2 \rceil^{m_0} e^{-u/96} \end{aligned} \quad (5)$$

for any $\kappa \in \{1, \dots, d\}^{m_0}$. Let E be the event that some hypothesis in \tilde{S} has error below $1/32$,

$$E = \{(\kappa, \lambda) \in \{1, \dots, d\}^{m_0+u}: \exists v, \mathbf{er}_{\tilde{P}, g}(\tilde{h}_{\kappa, v}) < 1/32\}.$$

(Notice that this event is independent of the examples $\lambda \in \{1, \dots, d\}^u$ that are used to assess the function in \tilde{S} .) For $\kappa \in \{1, \dots, d\}^{m_0}$ and $\lambda \in \{1, \dots, d\}^u$, let $\tilde{A}_{\kappa, \lambda, g}$ denote

$$\tilde{A}(\kappa_1, g(\kappa_1), \dots, \kappa_{m_0}, g(\kappa_{m_0}), \lambda_1, g(\lambda_1), \dots, \lambda_u, g(\lambda_u)).$$

Then (5) and the definition of u imply

$$\Pr(\mathbf{er}_{\tilde{P}, g}(\tilde{A}_{\kappa, \lambda, g}) > 1/8 | E) \leq 2\lceil 1/\alpha + 1/2 \rceil^{m_0} e^{-u/96} \leq 1/16, \quad (6)$$

where the probability is taken over all values of κ and λ conditioned on $(\kappa, \lambda) \in E$. But (4), which shows that $\Pr(\text{not } E) < 1/16$, and (6) imply

$$\begin{aligned} & \Pr(\mathbf{er}_{\tilde{P}, g}(\tilde{A}_{\kappa, \lambda, g}) > 1/8) \\ & \leq \Pr(\mathbf{er}_{\tilde{P}, g}(\tilde{A}_{\kappa, \lambda, g}) > 1/8 | E) + \Pr(\text{not } E) \\ & < 1/8. \end{aligned}$$

That is, \tilde{A} $(1/8, 1/8)$ -learns POWER_d using $m_0 + \lceil 96(\ln 32 + m_0 \ln \lceil 1/\alpha + 1/2 \rceil) \rceil$ examples, as claimed. Applying Theorem 8, this implies

$$\begin{aligned} m_0 + \lceil 96 \ln 32 + 96 m_0 \ln \lceil 1/\alpha + 1/2 \rceil \rceil &\geq d/2 \\ \Rightarrow m_0(1 + \lceil 96 \ln \lceil 1/\alpha + 1/2 \rceil \rceil) + 333 &\geq d/2 \\ \Rightarrow m_0 &\geq \frac{d/2 - 333}{2 + 96 \ln \lceil 1/\alpha + 1/2 \rceil}, \end{aligned}$$

which implies the lemma. \blacksquare

3.3. The Lower Bound

In this section, we combine Lemmas 5 and 10 to prove the following lower bound on the number of examples necessary for learning with observation noise. Obviously the constants have not been optimized.

THEOREM 11. *Suppose F is a set of $[0, 1]$ -valued functions defined on X , \mathcal{D} is an admissible noise distribution class with total variation function v , $0 < \gamma < 1$, $0 < \varepsilon \leq \gamma/65$, $0 < \delta \leq 1/32$, $\sigma \in \mathbb{R}^+$, and $d \in \mathbb{N}$. If $\text{fat}_F(\gamma) \geq d > 1000$, then any algorithm that $(\varepsilon, \delta, \sigma)$ -learns F with noise \mathcal{D} requires at least m_0 examples, where*

$$m_0 > \min \left\{ \frac{d}{1152 \ln(2 + dv(\sigma)/17)}, \frac{d}{1152 \ln(d/238)}, \frac{d}{576 \ln(35/\gamma)} \right\}. \quad (7)$$

In particular, if

$$v(\sigma) > \max(1/14, 101/(d\sqrt{\gamma})), \quad (8)$$

then

$$m_0 > \frac{d}{1152 \ln(2 + dv(\sigma)/17)}.$$

This theorem shows that if there is a $\gamma > 0$ such that $\text{fat}_F(\gamma)$ is infinite then we can choose ε , δ , and σ for which $(\varepsilon, \delta, \sigma)$ -learning is impossible from a finite sample. Similarly, if $\text{fat}_F(\gamma)$ grows faster than polynomially in $1/\gamma$, we can fix σ and Theorem 11 implies that the number of examples necessary for learning must grow faster than polynomially in $1/\varepsilon$. This proves the ‘‘only if’’ parts of the characterization theorem (Theorem 3).

We will use the following lemma.

LEMMA 12. *If $x, y, z > 0$, $\gamma z \geq 1$, $w \geq 1$, and $x > z/\ln(w(1 + xy))$, then $x > z/(2 \ln(w(1 + yz)))$.*

Proof. Suppose $x < z/(2 \ln(w(1 + yz)))$. Then, since $x \ln(w(1 + xy))$ is an increasing function of x , we have

$$\begin{aligned} x \ln(w(1 + xy)) &< \left(\frac{z}{2 \ln(w(1 + yz))} \right) \ln \left(w \left(1 + \frac{yz}{2 \ln(w(1 + yz))} \right) \right) \\ &= z \left(\frac{\ln(w(1 + (yz/(2 \ln(w(1 + yz))))))}{2 \ln(w(1 + yz))} \right). \end{aligned}$$

But $yz \geq 1$, so

$$w \left(1 + \frac{yz}{2 \ln(w(1 + yz))} \right) < w^2(1 + yz)^2,$$

which implies $x \ln(w(1 + xy)) < z$, a contradiction. \blacksquare

Proof (of Theorem 11). Set $\varepsilon = \gamma/65$ and $\delta = 1/32$. Suppose a learning algorithm can $(\varepsilon, \delta, \sigma)$ -learn F from m examples with noise \mathcal{D} . Lemma 5 shows that, provided

$$\alpha \leq \min(\delta/(v(\sigma)m), 2\varepsilon), \quad (9)$$

then there is a learning algorithm that can $(2\varepsilon, 2\delta)$ -learn $Q_\alpha(F)$ from m examples. From the definition of fat-shattering, $\text{fat}_F(\gamma) \geq d$ implies $\text{fat}_{Q_\alpha(F)}(\gamma - \alpha/2) \geq d$. Furthermore, since $\varepsilon = \gamma/65$, if Inequality (9) is satisfied, we have

$$(\gamma - \alpha/2)/32 \geq (\gamma - \gamma/65)/32 = 2\varepsilon.$$

Lemma 10 shows that, if an algorithm can $(2\varepsilon, 2\delta)$ -learn $Q_\alpha(F)$ from m examples (when $2\varepsilon \leq (\gamma - \alpha/2)/32$ and $2\delta \leq 1/16$), then

$$m \geq \frac{d - 666}{4 + 192 \ln \lceil 1/\alpha + 1/2 \rceil}. \quad (10)$$

That is, if Inequality (9) is satisfied, we must have m at least this large.

Using a case-by-case analysis, in each case choosing α to satisfy Inequality (9), we will show that m is larger than at least one of the terms in (7).

Consider the two cases $2\varepsilon \geq \delta/(v(\sigma)m)$ and $2\varepsilon < \delta/(v(\sigma)m)$.

Case 1. ($2\varepsilon \geq \delta/(v(\sigma)m)$). If we set $\alpha = \delta/(v(\sigma)m)$, Inequality (9) is satisfied, so

$$\begin{aligned} m &\geq \frac{d - 666}{4 + 192 \ln \lceil v(\sigma)m/\delta + 1/2 \rceil} \\ &> \frac{d/3}{4 + 192 \ln(32v(\sigma)m + 3/2)} \\ &= \frac{d}{12 + 576 \ln(\frac{3}{2}(1 + 64v(\sigma)m/3))}. \end{aligned} \quad (11)$$

Consider the two cases $v(\sigma) > 3/64$ and $v(\sigma) \leq 3/64$. First, suppose $v(\sigma) > 3/64$. Using Lemma 12 with $x = m$, $z = d/576$, $w = \frac{3}{2}e^{12/576}$, and $y = 64v(\sigma)/3$ (so $yz > d/576 > 1$), we have

$$\begin{aligned} m &> d/(1152 \ln(\frac{3}{2}e^{12/576} + e^{12/576} dv(\sigma)/18)) \\ &> d/(1152 \ln(2 + dv(\sigma)/17)), \end{aligned}$$

which is the first term in the minimum of Inequality (7). Now suppose that $v(\sigma) \leq 3/64$. Then (11) implies

$$m > d/(12 + 576 \ln(\frac{3}{2}(1 + m))).$$

Using Lemma 12 with $x = m$, $z = d/576$, $w = \frac{3}{2}e^{12/576}$, and $y = 1$ (and noting that $yz = d/576 > 1$), we have

$$\begin{aligned} m &> d/(1152 \ln(\frac{3}{2}e^{12/576}(1 + d/576))) \\ &> \frac{d}{1152 \ln(d/238)}, \end{aligned}$$

which is the second term in the minimum of Inequality (7).

Case 2. ($2\varepsilon < \delta/(v(\sigma)m)$). If we set $\alpha = 2\varepsilon$, Inequality (9) is satisfied, so Inequality (10) implies

$$\begin{aligned} m &> \frac{d - 666}{4 + 192 \ln \lceil 1/(2\varepsilon) + 1/2 \rceil} \\ &> \frac{d}{12 + 576 \ln(65/2\gamma + 3/2)} \\ &> \frac{d}{576 \ln(35/\gamma)}, \end{aligned}$$

which is the third term in the minimum of Inequality (7).

We now use Inequality (7) to prove the second part of the theorem. If

$$1152 \ln(2 + dv(\sigma)/17) > 1152 \ln(d/238) \quad (12)$$

and

$$1152 \ln(2 + dv(\sigma)/17) > 576 \ln(35/\gamma) \quad (13)$$

then

$$m_0 > \frac{d}{1152 \ln(2 + dv(\sigma)/17)}.$$

So it suffices to show that (12) and (13) are implied by (8). Indeed, we have that

$$\begin{aligned} v(\sigma) &> 1/14 \\ &\Rightarrow dv(\sigma)/17 > d/238 \\ &\Rightarrow 2 + dv(\sigma)/17 > d/238, \end{aligned}$$

which implies (12). Similarly,

$$\begin{aligned} v(\sigma) &> 101/(d\sqrt{\gamma}) \\ &\Rightarrow 2 + dv(\sigma)/17 > \sqrt{35/\gamma}, \end{aligned}$$

which implies (13). ■

4. UPPER BOUND

In this section, we prove an upper bound on the number of examples required for learning with observation noise, finishing the proof of Theorem 3. For $n \in \mathbb{N}$, $v, w \in \mathbb{R}^n$, let

$$d(v, w) = \frac{1}{n} \sum_{i=1}^n |v_i - w_i|.$$

For $U \subseteq \mathbb{R}^n$, $\varepsilon > 0$, we say $C \subseteq \mathbb{R}^n$ is an ε -cover of U if and only if for all $v \in U$, there exists $w \in C$ such that $d(v, w) \leq \varepsilon$, and we denote by $\mathcal{N}(\varepsilon, U)$ the size of the smallest ε -cover of U (the ε -covering number of U).

For a function $f: X \rightarrow [0, 1]$, define $l_f: X \times \mathbb{R} \rightarrow \mathbb{R}$ by $l_f(x, y) = (f(x) - y)^2$, and if $F \subseteq [0, 1]^X$, let $l_F = \{l_f: f \in F\}$.

If W is a set, $f: W \rightarrow \mathbb{R}$, and $w \in W^m$, let $f_{|w} \in \mathbb{R}^m$ denote $(f(w_1), \dots, f(w_m))$. Finally, if F is a set of functions from W to \mathbb{R} , let $F_{|w} \subseteq \mathbb{R}^m$ be defined by $F_{|w} = \{f_{|w}: f \in F\}$.

The following theorem is due to Haussler [15, Theorem 3, p. 107]; it is an improvement of a result of Pollard [22]. We say a function class is **PH-permissible** if it satisfies the mild measurability condition defined in Haussler's Section 9.2 in [15]. We say a class F of real-valued functions is **permissible** if the class l_F is PH-permissible. This implies that the class $l_F^a = \{(x, y) \mapsto |f(x) - y|: f \in F\}$ is PH-permissible, since the square root function on \mathbb{R}^+ is measurable.

THEOREM 13 [15]. *Let Y be a set and G a PH-permissible class of $[0, M]$ -valued functions defined on $Z = X \times Y$, where $M \in \mathbb{R}^+$. For any $\alpha > 0$ and any distribution P on Z ,*

$$\begin{aligned} P^m \left\{ z \in Z^m: \exists g \in G, \left| \frac{1}{m} \sum_{i=1}^m g(z_i) - \int_Z g dP \right| > \alpha \right\} \\ \leq 4 \max_{z \in Z^{2m}} (\mathcal{N}(\alpha/16, G_{|z})) e^{-\alpha^2 m / (64M^2)}. \end{aligned}$$

COROLLARY 14. *Let F be a permissible class of $[0, 1]$ -valued functions defined on X . Let $Y = [a, b]$ with $a \leq 0$ and $b \geq 1$, and let $Z = X \times Y$. There is a mapping B from $(0, 1) \times \bigcup_i Z^i$ to $[0, 1]^X$ such that, for any $0 < \varepsilon < 1$ and any distribution P on Z ,*

$$P^m \left\{ z \in Z^m: \int_Z l_{B(\varepsilon, z)} dP \geq \inf_{f \in F} \int_Z l_f dP + \varepsilon \right\} \leq 4 \max_{z \in Z^{2m}} (\mathcal{N}(\varepsilon/48, l_{F|z})) e^{-\varepsilon^2 m / (576(b-a)^4)}.$$

The proof is similar to the proof of Haussler's Lemma 1 in [15].

Proof. For a sequence $z = (z_1, \dots, z_m)$, let the mapping B return a function f^* from F that satisfies

$$\frac{1}{m} \sum_{i=1}^m l_{f^*}(z_i) < \inf_{f \in F} \frac{1}{m} \sum_{i=1}^m l_f(z_i) + \varepsilon/3. \quad (14)$$

Let $M = (b-a)^2$. Theorem 13 implies that, with probability at least

$$1 - 4 \max \mathcal{N}(\varepsilon/48, l_{F|z}) e^{-\varepsilon^2 m / (576M^2)},$$

we have

$$\left| \frac{1}{m} \sum_{i=1}^m l_{f^*}(z_i) - \int_Z l_{f^*} dP \right| < \varepsilon/3 \quad (15)$$

and

$$\left| \inf_{f \in F} \frac{1}{m} \sum_{i=1}^m l_f(z_i) - \inf_{f \in F} \int_Z l_f dP \right| < \varepsilon/3. \quad (16)$$

By the triangle inequality for absolute difference on the reals, (14), (15), and (16) imply

$$\left| \int_Z l_{f^*} dP - \inf_{f \in F} \int_Z l_f dP \right| < \varepsilon. \quad \blacksquare$$

The following result follows trivially from Alon, Ben-David, Cesa-Bianchi, and Haussler's Lemmas 14 and 15 [1].

THEOREM 15 [1]. *If F is a class of $[0, 1]$ -valued functions defined on X , $0 < \varepsilon < 1$, and $m \in \mathbb{N}$, then for all x in X^m ,*

$$\mathcal{N}(\varepsilon, F|_x) \leq 2(mb^2)^{\log c},$$

where $b = \lceil 2/\varepsilon \rceil + 1$ and

$$c = \sum_{i=1}^{\text{fat}_F(\varepsilon/4)} \binom{m}{i} b^i.$$

COROLLARY 16. *For F defined as in Theorem 15, if $0 < \varepsilon < 1/2$ and $m \geq \text{fat}_F(\varepsilon/4)/2$, then for all x in X^m*

$$\mathcal{N}(\varepsilon, F|_x) \leq \exp \left(\frac{2}{\ln 2} \text{fat}_F(\varepsilon/4) \ln^2 \frac{9m}{\varepsilon^2} \right).$$

Proof. Let $d = \text{fat}_F(\varepsilon/4)$. If $d = 0$ then any f_1 and f_2 in F have $|f_1(x) - f_2(x)| < \varepsilon/2$, so $\mathcal{N}(\varepsilon, F|_x) \leq 1$ in this case. Assume then that $d \geq 1$. We have $b < 3/\varepsilon$ and

$$\begin{aligned} \log c &< \log \sum_{i=1}^d \binom{m}{i} (3/\varepsilon)^i \\ &< \log \left(d \binom{m}{d} (3/\varepsilon)^d \right) \\ &< \log(d(3m/\varepsilon)^d) \\ &< d \log(3m/\varepsilon) + \log d. \end{aligned}$$

So we have

$$\begin{aligned} \ln \mathcal{N}(\varepsilon, F|_x) &\leq \ln 2 + (d \log(3m/\varepsilon) + \log d) \ln(9m/\varepsilon^2) \\ &< 2d \ln(3m/\varepsilon) \ln(9m/\varepsilon^2) / \ln 2 \\ &< 2d \ln^2(9m/\varepsilon^2) / \ln 2. \quad \blacksquare \end{aligned}$$

Note that the bound of Corollary 14 involves covering numbers of l_F , whereas Corollary 16 bounds covering numbers of F . This was handled in [1] in the case of probabilistic concepts (where the $Y = [a, b]$ in Corollary 14 is replaced by $Y = \{0, 1\}$) by showing that in that case, $\text{fat}_{l_F}(\gamma) \leq \text{fat}_F(\gamma/2)$. In the following lemma, we relate the covering numbers of l_F and of F .⁵

LEMMA 17. *Choose a set F of functions from X to $[0, 1]$. Then for any $\varepsilon > 0$, for any $m \in \mathbb{N}$, if $a \leq 0$ and $b \geq 1$,*

$$\max_{z \in (X \times [a, b])^m} \mathcal{N}(\varepsilon, (l_F)|_z) \leq \max_{x \in X^m} \mathcal{N} \left(\frac{\varepsilon}{3|b-a|}, F|_x \right).$$

Proof. We show that, for any sequence z of (x, y) pairs in $X \times [a, b]$ and any functions f and g , if the restrictions of f and g to x are close, then the restrictions of l_f and l_g to z are close. Thus, given a cover of $F|_x$, we can construct a cover of $l_{F|z}$ that is no bigger.

Now, choose $(x_1, y_1), \dots, (x_m, y_m) \in X \times [a, b]$, and $f, g: X \rightarrow [0, 1]$. We have

⁵ Recently, Gurvits and Koiran have proved a result relating the fat-shattering functions of l_F and F [14].

$$\begin{aligned} & \frac{1}{m} \sum_{i=1}^m |(g(x_i) - y_i)^2 - (f(x_i) - y_i)^2| \\ &= \frac{1}{m} \sum_{i=1}^m |(g(x_i) - y_i)^2 - ((f(x_i) - g(x_i)) + g(x_i) - y_i)^2| \\ &= \frac{1}{m} \sum_{i=1}^m |(f(x_i) - g(x_i))^2 - 2(f(x_i) - g(x_i))(g(x_i) - y_i)| \\ &\leq \frac{1}{m} \sum_{i=1}^m ((f(x_i) - g(x_i))^2 + 2|f(x_i) - g(x_i)| |g(x_i) - y_i|) \\ &\leq \frac{1}{m} \sum_{i=1}^m 3|b - a| |f(x_i) - g(x_i)|. \end{aligned}$$

Thus if $x = (x_1, \dots, x_m) \in X^m$, $z = (x_1, y_1, \dots, x_m, y_m) \in (X \times [a, b])^m$, and $d(f|_x, g|_x) \leq \varepsilon/(3|b - a|)$, then $d(l_{f|z}, l_{g|z}) \leq \varepsilon$. So if S is an $\varepsilon/(3|b - a|)$ -cover of $F|_x$, we can construct an ε -cover T of $l_{F|z}$ as

$$T = \{((u_1 - y_1)^2, \dots, (u_m - y_m)^2) : u \in S\}.$$

Since $(x_1, y_1), \dots, (x_m, y_m)$ was chosen arbitrarily, this completes the proof. ■

In our proof of upper bounds on the number of examples needed for learning, we will make use of the following lemma.

LEMMA 18. For any $y_1, y_2, y_4, \delta > 0$, and $y_3 \geq 1$, if

$$m \geq \frac{2}{y_4} \left(4y_2 \left(4 + \ln \left(\frac{y_2 y_3}{y_4} \right) \right)^2 + \ln \frac{y_1}{\delta} \right),$$

then

$$y_1 \exp(y_2 \ln^2(y_3 m) - y_4 m) \leq \delta.$$

Proof. The assumed lower bound on m implies that

$$m \geq \frac{2}{y_4} \ln \frac{y_1}{\delta} \tag{17}$$

and

$$m \geq \frac{8y_2}{y_4} \left(2 \ln(4\sqrt{2}) + \ln \left(\frac{y_2 y_3}{y_4} \right) \right)^2.$$

Taking square roots of the latter inequality and fiddling a little with the second term, we get

$$\sqrt{m} \geq 2\sqrt{2} \sqrt{y_2/y_4} (2 \ln(4\sqrt{2} \sqrt{y_2/y_4}) + \ln y_3).$$

Setting $b = (1/(4\sqrt{2})) \sqrt{y_4/y_2}$, the previous inequality implies that

$$\sqrt{m}(1 - 2\sqrt{2y_2/y_4}b) \geq \sqrt{2} \sqrt{y_2/y_4} (2 \ln(1/b) + \ln y_3)$$

which trivially yields

$$\sqrt{m} \geq \sqrt{2} \sqrt{y_2/y_4} (2(b\sqrt{m} + \ln(1/b)) + \ln y_3).$$

The above inequality, using the fact [24] that for all $a, b > 0$, $\ln a \leq ab + \ln(1/b)$, implies that

$$\begin{aligned} \sqrt{m} &\geq \sqrt{2} \sqrt{y_2/y_4} (2 \ln \sqrt{m} + \ln y_3) \\ &= \sqrt{2} \sqrt{y_2/y_4} \ln(y_3 m). \end{aligned}$$

Squaring both sides and combining with (17), we get

$$m \geq \frac{1}{y_4} (y_2 \ln^2(y_3 m) + \ln(y_1/\delta)).$$

Solving for δ completes the proof. ■

We can now present the upper bound. Again, the constants have not been optimized.

THEOREM 19. For any permissible class F of functions from X to $[0, 1]$, there is a learning algorithm A such that, for all bounded admissible distribution classes \mathcal{D} with support function s , for all probability distributions P on X , and for all $0 < \varepsilon < 1/2$, $0 < \delta < 1$, and $\sigma > 0$, if $d = \text{fat}_F(\varepsilon^2/(576(s(\sigma) + 1)))$, then $A(\varepsilon, \delta, \sigma)$ -learns F from

$$\frac{1152(1 + s(\sigma))^4}{\varepsilon^4} \left(12d \left(25 + \ln \frac{d(1 + s(\sigma))^6}{\varepsilon^8} \right)^2 + \ln \frac{4}{\delta} \right)$$

examples with noise \mathcal{D} .

Proof. Let B be the mapping from Corollary 14. Choose $0 < \varepsilon < 1/2$, $0 < \delta < 1$, and $\sigma > 0$. Let $\varepsilon_0 = \varepsilon^2$. Let D be a distribution in \mathcal{D} with variance σ^2 and support contained in $[c, d]$, so $d - c \leq s(\sigma)$. Choose a distribution P on X and a function $f \in F$.

For $x \in X^m$ and $\eta \in [c, d]^m$, let $B_{x,\eta} = B(\varepsilon_0, \text{sam}(x, \eta, f))$. Define the event

BAD

$$= \left\{ (x, \eta) \in (X^m \times [c, d]^m) : \right.$$

$$\left. \int_X \int_{[c, d]} [B_{x,\eta}(u) - (f(u) + \kappa)]^2 dD(\kappa) dP(u) \geq \sigma^2 + \varepsilon_0 \right\}.$$

Since D has variance σ^2 and mean 0,

$$\inf_{g \in F} \int_X \int_{[c, d]} (g(u) - (f(u) + \kappa))^2 dD(\kappa) dP(u) = \sigma^2,$$

so

BAD

$$\begin{aligned} &= \left\{ (x, \eta): \int_X \int_{[c, d]} [B_{x, \eta}(u) - (f(u) + \kappa)]^2 dD(\kappa) dP(u) \right. \\ &\geq \left. \inf_{g \in F} \int_X \int_{[c, d]} [g(u) - (f(u) + \kappa)]^2 dD(\kappa) dP(u) + \varepsilon_0 \right\}. \end{aligned}$$

The random variable $f(u) + \kappa$ has a distribution on $[c, 1 + d]$, determined by the distributions P and D and the function f . Thus, by Corollary 14,

$$\begin{aligned} \Pr(\text{BAD}) &\leq 4 \left(\max_{z \in (X \times [c, 1 + d])^{2m}} \mathcal{N}(\varepsilon_0/48, (I_F)_|z) \right) \\ &\quad \times \exp\left(\frac{-\varepsilon_0^2 m}{576(1 + s(\sigma))^4}\right). \end{aligned}$$

Lemma 17 implies

$$\begin{aligned} \Pr(\text{BAD}) &\leq 4 \left(\max_{x \in X^{2m}} \mathcal{N}\left(\frac{\varepsilon_0}{144(1 + s(\sigma))}, F_{|x}\right) \right) \\ &\quad \times \exp\left(\frac{-\varepsilon_0^2 m}{576(1 + s(\sigma))^4}\right). \end{aligned}$$

Applying Corollary 16, if

$$d = \text{fat}_F\left(\frac{\varepsilon_0}{576(1 + s(\sigma))}\right),$$

and $m \geq d/2$, then

$\Pr(\text{BAD})$

$$\begin{aligned} &\leq 4 \exp\left(\frac{2}{\ln 2} d \ln^2 \frac{373248m(1 + s(\sigma))^2}{\varepsilon_0^2} \right. \\ &\quad \left. - \frac{\varepsilon_0^2 m}{576(1 + s(\sigma))^4}\right). \end{aligned} \quad (18)$$

For any particular $x \in X^m$, $\eta \in [c, d]^m$,

$$\begin{aligned} &\int_X \int_{[c, d]} (B_{x, \eta}(u) - (f(u) + \kappa))^2 dD(\kappa) dP(u) \\ &= \int_X \int_{[c, d]} (B_{x, \eta}(u) - f(u))^2 dD(\kappa) dP(u) \\ &\quad - 2 \int_X \int_{[c, d]} (B_{x, \eta}(u) - f(u)) \kappa dD(\kappa) dP(u) + \sigma^2 \\ &= \int_X \int_{[c, d]} (B_{x, \eta}(u) - f(u))^2 dD(\kappa) dP(u) + \sigma^2 \end{aligned}$$

because of the independence of the noise and the fact that it has zero mean. Thus

$$\begin{aligned} \text{BAD} &= \left\{ (x, \eta) \in (X^m \times [a, b]^m): \right. \\ &\quad \left. \int_X [B_{x, \eta}(u) - f(u)]^2 dP(u) \geq \varepsilon_0 \right\}. \end{aligned}$$

If

$$m \geq \frac{1152(1 + s(\sigma))^4}{\varepsilon_0^2} \left(12d \left(25 + \ln \frac{d(1 + s(\sigma))^6}{\varepsilon_0^4} \right)^2 + \ln \frac{4}{\delta} \right), \quad (19)$$

then applying Lemma 18, with $y_1 = 4$, $y_2 = 2d/\ln 2$, $y_3 = 373248(1 + s(\sigma))^2/\varepsilon_0^2$, and $y_4 = \varepsilon_0^2/(576(1 + s(\sigma))^4)$, we have that (18) and (19) imply

$$P^m \times D^m \left\{ (x, \eta): \int_X (B_{x, \eta}(u) - f(u))^2 dP(u) \geq \varepsilon_0 \right\} < \delta. \quad (20)$$

From Jensen's inequality,

$$\begin{aligned} &\left\{ (x, \eta): \int_X |B_{x, \eta}(u) - f(u)| dP(u) \geq \sqrt{\varepsilon_0} \right\} \\ &\subseteq \left\{ (x, \eta): \int_X (B_{x, \eta}(u) - f(u))^2 dP(u) \geq \varepsilon_0 \right\}, \end{aligned}$$

so if $m \geq m_0(\varepsilon, \delta, \sigma)$,

$$P^m \times D^m \left\{ (x, \eta): \int_X |(B(\varepsilon^2, \text{sam}(x, \eta, f)))(u) - f(u)| dP(u) \geq \varepsilon \right\} < \delta,$$

where

$$m_0(\varepsilon, \delta, \sigma) = \frac{1152(1 + s(\sigma))^4}{\varepsilon^4} \times \left(12d \left(25 + \ln \frac{d(1 + s(\sigma))^6}{\varepsilon^8} \right)^2 + \ln \frac{4}{\delta} \right),$$

and

$$d = \text{fat}_F \left(\frac{\varepsilon^2}{576(1 + s(\sigma))} \right).$$

Now, let A be the algorithm that counts the number m of examples it receives and chooses ε_1 such that $m_0(\varepsilon_1, 1, 0) = m$. This is always possible, since d and, hence, m_0 are nonincreasing functions of ε . Algorithm A then passes ε_1^2 and the examples to the mapping B , and returns B 's hypothesis. Since $s(\sigma)$ is a nondecreasing function of σ , m_0 is a nondecreasing function of $1/\varepsilon$, $1/\delta$, and σ , so for any ε , δ , and σ satisfying $m_0(\varepsilon, \delta, \sigma) \leq m$, we must have $\varepsilon \geq \varepsilon_1$. It follows that, for any ε , δ , and σ for which A sees at least $m_0(\varepsilon, \delta, \sigma)$ examples, if P is a distribution on $X \times Y$ and $D \in \mathcal{D}$ has variance σ^2 then

$$P^m \times D^m \left\{ (x, \eta) : \int_X |(A(\text{sam}(x, \eta, f)))(u) - f(u)| dP(u) \geq \varepsilon \right\} < \delta,$$

completing the proof. ■

As an immediate consequence of Theorem 19, if F has a finite fat-shattering function and \mathcal{D} is a bounded admissible distribution class, then F is learnable with observation noise \mathcal{D} . The following corollary provides the one implication in Theorem 3 we have yet to prove.

COROLLARY 20. *Let F be a class of functions from X to $[0, 1]$. Let p be a polynomial, and suppose $\text{fat}_F(\gamma) < p(1/\gamma)$ for all $0 < \gamma < 1$. Then for any almost-bounded admissible noise distribution class \mathcal{D} , F is small-sample learnable with noise \mathcal{D} .*

Proof. We will show that Algorithm A from Theorem 19 can $(\varepsilon, \delta, \sigma)$ -learn F from a polynomial number of examples with noise \mathcal{D} .

Let $s: \mathbb{R}^+ \rightarrow \mathbb{R}^+$ (we will define s later). Choose $0 < \varepsilon, \delta < 1, \sigma > 0$. Fix a distribution P on X , a function f in F , and a noise distribution D in \mathcal{D} with variance σ^2 .

Construct a distribution D_s from D as follows. Let ϕ be the pdf of D . Define the pdf ϕ_s of D_s as

$$\phi_s(x) = \begin{cases} \frac{\phi(x)}{\int_{-s(\sigma)/2}^{s(\sigma)/2} \phi(x) dx} & \text{if } -s(\sigma)/2 < x < s(\sigma)/2 \\ 0 & \text{otherwise.} \end{cases}$$

Since \mathcal{D} is an almost-bounded admissible class, there are universal constants $s_0, c_0 \in \mathbb{R}^+$ such that, if $s(\sigma) > s_0\sigma$,

$$\int_{-s(\sigma)/2}^{s(\sigma)/2} \phi(x) dx \geq 1 - c_0 e^{-s(\sigma)/\sigma}.$$

Let $I = \int_{-s(\sigma)/2}^{s(\sigma)/2} \phi(x) dx$. The total variation distance between D and D_s is

$$\begin{aligned} d_{\text{TV}}(D, D_s) &= \int_{-\infty}^{\infty} |\phi(x) - \phi_s(x)| dx \\ &= 1 - I + \int_{-s(\sigma)/2}^{s(\sigma)/2} |\phi(x) - \phi_s(x)| dx \\ &= 1 - I + |1 - I| \int_{-s(\sigma)/2}^{s(\sigma)/2} \phi(x) dx \\ &= 2(1 - I) \\ &\leq 2c_0 e^{-s(\sigma)/\sigma}. \end{aligned} \tag{21}$$

For some m in \mathbb{N} , fix $x \in X^m$ and define the event

$$E_1 = \{ \eta \in \mathbb{R}^m : \mathbf{er}_{P,f}(A(\text{sam}(x, \eta, f))) \geq \varepsilon \}.$$

Then (21) and Lemma 7 show that

$$D^m(E_1) \leq D_s^m(E_1) + mc_0 \exp(-s(\sigma)/\sigma).$$

If we choose $s(\sigma) = \sigma(s_0 + |\ln(mc_0/\delta)|)$, then (21) holds and $s(\sigma) \geq \sigma \ln(2mc_0/\delta)$, so $D^m(E_1) \leq D_s^m(E_1) + \delta/2$. Since this is true for any $x \in X^m$,

$$P^m \times D^m(E_2) \leq P^m \times D_s^m(E_2) + \delta/2,$$

where

$$E_2 = \{ (x, \eta) \in X^m \times \mathbb{R}^m : \mathbf{er}_{P,f}(A(\text{sam}(x, \eta, f))) \geq \varepsilon \}.$$

Clearly, D_s has mean 0, finite variance, and support contained in an interval of length $s(\sigma)$. From the proof of Theorem 19, there is a polynomial p_1 such that if

$$m \geq p_1(s(\sigma), d, 1/\varepsilon, \ln 1/\delta)$$

then

$$P^m \times D_s^m(E_2) < \delta/2. \tag{22}$$

Now, $\text{fat}_F(\gamma) < p(1/\gamma)$, so for some polynomial p_2 , $m > p_2(\sigma, 1/\varepsilon, \log(1/\delta), \log m)$ implies (22). Clearly, for some polynomial p_3 , if $m > p_3(\sigma, 1/\varepsilon, \log(1/\delta))$ then $P^m \times D^m(E_2) < \delta$. Since this is true for any P and any D in \mathcal{D} with variance σ^2 , Algorithm $A(\varepsilon, \delta, \sigma)$ -learns F with noise \mathcal{D} from $p_3(\sigma, 1/\varepsilon, \log(1/\delta))$ examples. ■

5. AGNOSTIC LEARNING

In this section, we consider an agnostic learning model, a model of learning in which assumptions about the target function and observation noise are removed. In this model, we assume labelled examples (x, y) are generated by some joint distribution P on $X \times [0, 1]$. The agnostic learning problem can be viewed as the problem of learning a real-valued function f with observation noise when the constraints on the noise are relaxed—in particular, we no longer have the constraint that the noise is independent of the value $f(x)$. This model has been studied in [15, 17].

If h is a $[0, 1]$ -valued function defined on X , define the **error of h with respect to P** as

$$\text{er}_P(h) = \int_{X \times [0, 1]} |h(x) - y| dP(x, y).$$

We require that the learner chooses a function with error little worse than the best function in some “touchstone” function class F . Notice that the learner is not restricted to choose a function from F ; the class F serves only to provide a performance measurement standard (see [17]).

DEFINITION 21. Suppose F is a class of $[0, 1]$ -valued functions defined on X , P is a probability distribution on $X \times [0, 1]$, $0 < \varepsilon, \delta < 1$, and $m \in \mathbb{N}$. We say a learning algorithm $L = (A, D_Z)$ **(ε, δ) -learns in the agnostic sense with respect to F from m examples** if, for all distributions P on $X \times [0, 1]$,

$$(P^m \times D_Z^m)\{(x, y, z) \in X^m \times [0, 1]^m \times Z^m: \text{er}_P(A(x, y, z)) \geq \inf_{f \in F} \text{er}_P(f) + \varepsilon\} < \delta.$$

The function class F is **agnostically learnable** if there is a learning algorithm L and a function $m_0: (0, 1) \times (0, 1) \rightarrow \mathbb{N}$ such that, for all $0 < \varepsilon, \delta < 1$, algorithm $L(\varepsilon, \delta)$ -learns in the agnostic sense with respect to F from $m_0(\varepsilon, \delta)$ examples. If, in addition, m_0 is bounded by a polynomial in $1/\varepsilon$ and $1/\delta$, we say that F is **small-sample agnostically learnable**.

The following result is analogous to the characterization theorem of Section 2.

THEOREM 22. *Suppose F is a permissible class of $[0, 1]$ -valued functions defined on X . Then F is agnostically learnable if and only if its fat-shattering function is finite, and F is small-sample agnostically learnable if and only if there is a polynomial p such that $\text{fat}_F(\gamma) < p(1/\gamma)$ for all $\gamma > 0$.*

Alon *et al.*'s proof in [1] that finiteness of the fat-shattering function of the class I_F is sufficient for learnability of a class F of probabilistic concepts also shows that this condition is sufficient for the agnostic learnability of a class F of real-valued functions. A simpler version of Lemma 17 then shows that finiteness of the fat-shattering function of F suffices for agnostic learnability.

If the “loss” of the learning algorithm was measured with $(h(x) - y)^2$, instead of $|h(x) - y|$, then the necessity part of Theorem 22 would follow from the results of Kearns and Schapire [16].

The following result proves the “only if” parts of the theorem.

THEOREM 23. *Let F be a class of $[0, 1]$ -valued functions defined on X . Suppose $0 < \gamma < 1$, $0 < \varepsilon \leq \gamma/65$, $0 < \delta \leq 1/16$, and $d \in \mathbb{N}$. If $\text{fat}_F(\gamma) \geq d > 1000$, then any learning algorithm that (ε, δ) -learns in the agnostic sense with respect to F requires at least m_0 examples, where*

$$m_0 > \frac{d}{576 \ln(35/\gamma)}.$$

Proof. The proof is similar to, although simpler than, the argument in Section 3. We will show that the agnostic learning problem is not much harder than the problem of learning a quantized version of the function class F and then apply Lemma 10.

Set $\varepsilon = \gamma/65$ and $\delta = 1/16$. Consider the class of distributions P on $X \times [0, 1]$ for which there exists an f in F such that, for all $x \in X$,

$$P(y|x) = \begin{cases} 1 & y = Q_{2\varepsilon}(f(x)) \\ 0 & \text{otherwise.} \end{cases}$$

Fix a distribution P in this class. Let L be a randomized learning algorithm that can (ε, δ) -learn in the agnostic sense with respect to F . Then

$$\Pr(\text{er}_P(L) \geq \inf_{f \in F} (\text{er}_P(f)) + \varepsilon) < \delta,$$

where $\text{er}_P(L)$ is the error of the function that the learning algorithm chooses. But the definition of P ensures that $\inf_{f \in F} \text{er}_P(f) \leq \varepsilon$, so

$$\Pr(\text{er}_P(L) \geq 2\varepsilon) < \delta.$$

Since this is true for any distribution P that can be expressed as the product of a distribution on X and the generalized derivative of the indicator function on $[0, 1]$ of a function in $\mathcal{Q}_{2\varepsilon}(F)$, the learning algorithm L can $(2\varepsilon, \delta)$ -learn the quantized function class $\mathcal{Q}_{2\varepsilon}(F)$.

By hypothesis, $\text{fat}_F(\gamma) \geq d$, but then the definition of fat-shattering implies that $\text{fat}_{\mathcal{Q}_{2\varepsilon}(F)}(\gamma - \varepsilon) \geq d$. Since $\varepsilon = \gamma/65$, $2\varepsilon \leq (\gamma - \varepsilon)/32$. Also, $\delta = 1/16$, so Lemma 10 implies

$$\begin{aligned} m_0 &> \frac{d - 666}{4 + 192 \ln \lceil 1/(2\varepsilon) + 1/2 \rceil} \\ &> \frac{d}{12 + 576 \ln(65/(2\gamma) + 3/2)} \\ &> \frac{d}{576 \ln(35/\gamma)}. \quad \blacksquare \end{aligned}$$

With minor modifications, the proof of Theorem 19 yields the following analogous result for agnostic learning.

THEOREM 24. *Choose a permissible set F of functions from X to $[0, 1]$. There exists an algorithm A such that, for all $0 < \varepsilon < 1/2$, for all $0 < \delta < 1$, if $\text{fat}_F(\varepsilon/192) = d$, then A agnostically (ε, δ) -learns F from*

$$\frac{1152}{\varepsilon^2} \left(12d \left(23 + \ln \frac{d}{\varepsilon^4} \right)^2 + \ln \frac{4}{\delta} \right)$$

examples.

Proof Sketch. First, the analog of Corollary 14, where the expected absolute error is used to measure the “quality” of a hypothesis in place of the expected squared error, $b = 1$, and $a = 0$, can be proved using essentially the same argument. Second, the analog of Lemma 17, where l_F is replaced with a corresponding class constructed from absolute loss in place of l , $a = 0$, $b = 1$, and the $\varepsilon/(3|b - a|)$ of the upper bound is replaced with ε , also is obtained using a simpler, but similar, proof. These results are combined with Corollary 16 and Lemma 18 in much the same way as was done for Theorem 19. \blacksquare

6. DISCUSSION

All of our results can be extended easily to the case of $[L, U]$ -valued functions by scaling the parameters ε , γ , and σ to convert the learning problem to an equivalent $[0, 1]$ -valued learning problem.

It would be worthwhile to extend the characterization of learnability in terms of finiteness of the fat-shattering function to weaker noise models. It seems likely that it could be extended to the case of unbounded noise; perhaps the techniques used in [13] to prove uniform convergence with unbounded noise could be useful here.

There are several ways in which our results could be improved. The sample complexity upper bound in Theorem 19 increases at least as $1/\varepsilon^4$. It seems plausible that this rate is excessive; perhaps it is an artifact of the use of Jensen’s inequality in the proof. Obviously, the constants in our bounds are large. Another weakness of our bounds is the gap between constant factors in the argument of the fat-shattering function. If the domain X is infinite, this gap alone can lead to an arbitrarily large gap in the sample complexity bounds. Recent results [9] for agnostic learning narrow this gap to a factor of two.

The lower bound on the sample complexity of real-valued learning (Theorem 11) does not increase with $1/\varepsilon$ and $1/\delta$. In fact, the lower bound of that theorem is trivially true if the standard deviation of the noise is sufficiently small,⁶ i.e.,

$$\frac{1}{v(\sigma)} < de^{-d/1152}/17.$$

However, the following example shows that a condition of this form is essential and that when the noise variance is small there need be no dependence of the lower bound on the desired accuracy and confidence.

EXAMPLE. Fix $d \in \mathbb{N}$. Let the measurable sets S_j , $j = 0, \dots, d - 1$, form a partition of X (that is, $\bigcup_j S_j = X$, and $S_j \cap S_k = \emptyset$ if $j \neq k$). Consider the function class

$$F_d = \{f_{b_0, \dots, b_{d-1}} : b_i \in \{0, 1\}, i = 0, \dots, d - 1\}$$

of functions defined by

$$f_{b_0, \dots, b_{d-1}}(x) = \frac{3}{4} \sum_{j=0}^{d-1} 1_{S_j}(x) b_j + \frac{1}{8} \sum_{k=0}^{d-1} b_k 2^{-k},$$

where 1_{S_j} is the indicator function for S_j ($1_{S_j}(x) = 1$ iff $x \in S_j$). That is, the labels b_j determine the two most significant bits of the value of the function in S_j , and the d least significant bits of its value at any $x \in X$ encode the identity of the function. Clearly, for any $\gamma \leq 1/4$, $\text{fat}_{F_d}(\gamma) = d$.

With no observation noise, one example (x, y) suffices to learn F_d exactly, because the learning algorithm can identify the function from the d least significant bits of y . (As an aside, the union of these function classes, $F = \bigcup_{d=1}^{\infty} F_d$, has $\text{fat}_F(\gamma) = \infty$ for $\gamma \leq 1/4$, but any f in F can be identified from a single example (x, y) with no observation noise.⁷) One

⁶ Note that as the standard deviation gets small, the total variation of the density function must get large.

⁷ Thanks to David Haussler for suggesting this function class.

example also suffices with uniform observation noise, provided the variance is sufficiently small; if

$$\sigma < \frac{1}{2^{d+3} \sqrt{3}},$$

a learning algorithm that sees one example (x, y) and chooses the integral multiple of 2^{-d-2} that is closest to y will be able to identify the target function. That is, if

$$\frac{1}{v(\sigma)} < \frac{1}{2^{d+2} \sqrt{3}},$$

then $(\varepsilon, \delta, \sigma)$ -learning with uniform noise is possible from a single example, for any $\varepsilon, \delta \geq 0$.

Suppose the observation noise is gaussian, of variance σ^2 , and

$$\sigma < \frac{1}{2^{d+5/2} \sqrt{\log 4}}.$$

Consider the following algorithm. For each example (x, y) , the algorithm chooses the integral multiple of 2^{-d-2} that is closest to y and stores the corresponding function label (the d least significant bits). After m examples, it outputs the function with the most common label. The bound on σ and Inequality (1) (the bound on the area under the tails of the gaussian density) imply that, with probability at least $3/4$ a noisy observation is closer to the value $f(x)$ than to any other integral multiple of 2^{-d-2} . From Chernov bounds (see Theorem 9), if $m \geq 12 \log(1/\delta)$ the probability that the algorithm will store the correct label for fewer than half of the examples is less than δ . So this algorithm can $(\varepsilon, \delta, \sigma)$ -learn from $12 \log(1/\delta)$ examples, for any $\varepsilon \geq 0$. ■

The above example shows that a gap in the growth of the upper and lower bounds with $1/\varepsilon$ and $1/\delta$ is essential. However, the gap is unnecessarily large; a recent result relating several scale sensitive dimensions (Lemma 9 in [3]) implies improved lower bounds on the sample complexity of learning quantized function classes. In turn, these imply an improved general lower bound (of $\Omega(d/(\varepsilon \log^2(d/\varepsilon)))$) on the sample complexity of learning with observation noise that is valid if the noise variance is sufficiently large.

The example also shows that finiteness of the fat-shattering function is not necessary for learning real-valued functions without noise. However, the function classes that provide this counterexample are unnatural. We can interpret Theorem 3 as showing that if we change the definition of learning by requiring the learning algorithm to cope with additive observation noise, this rules out these unnatural

function classes. Similarly, the main result in [3] shows that, if the learning algorithm is constrained to return a function from the class that approximately interpolates the training examples, finiteness of the fat-shattering function is again necessary and sufficient for learning.

Simon [25] shows that a stronger notion of shattering provides a lower bound for the problem of learning without noise. However, the finiteness of this strong-fat-shattering function is not necessary for learnability, as the following example shows.

EXAMPLE. We say that a sequence x_1, \dots, x_d is **strongly γ -shattered** by F if there exist $u, l \in [0, 1]^d$ such that for each $b \in \{0, 1\}^d$ there is an $f \in F$ such that for each i , $u_i - l_i \geq 2\gamma$ and

$$f(x_i) = \begin{cases} u_i & \text{if } b_i = 1 \\ l_i & \text{if } b_i = 0. \end{cases}$$

For each γ , let

$$\text{sfat}_F(\gamma) = \max\{d \in \mathbb{N} : \exists x_1, \dots, x_d, F \text{ strongly } \gamma\text{-shatters } x_1, \dots, x_d\}$$

if such a maximum exists, and ∞ otherwise. If $\text{sfat}_F(\gamma)$ is finite for all γ , we say F has a **finite strong-fat-shattering function**.

Suppose $X = \mathbb{N}$. For each $q: \mathbb{N} \rightarrow \{0, 1\}$, let y_q be the element of $[0, 1]$ whose representation as a binary fraction is given by q , i.e., let $y_q = \sum_{i=1}^{\infty} q(i) 2^{-i}$. Also, let $f_q: \mathbb{N} \rightarrow [0, 1]$ be defined by

$$f_q(j) = \begin{cases} 3/4 + y_q/4 & \text{if } q(j) = 1 \\ 1/4 - y_q/4 & \text{if } q(j) = 0. \end{cases}$$

Let

$$Q = \{q \in \bigcup_{i=1}^{\infty} \{0, 1\}^i : \forall j_0 \exists j > j_0, q(j) = 0\}.$$

Informally, Q represents the set of all infinite binary sequences that do not end with repeating 1's. Each real number in $[0, 1)$ has a unique representation in Q [23]. Suppose $F = \{f_q : q \in Q\}$. Since X is countable, F is permissible. Trivially, $\text{fat}_F(1/4) = \infty$, so F is not learnable in any sense described in this paper. However, since for any $q_1, q_2 \in Q$ for which $q_1 \neq q_2$, for any $j \in \mathbb{N}$, $f_{q_1}(j) \neq f_{q_2}(j)$, trivially, $\text{sfat}_F(\gamma) = 1$ for all $\gamma < 1$, so neither the finiteness nor the polynomial growth of sfat_F characterizes learnability in any of the senses of this paper. ■

Simon provides examples in his paper that show that his general lower bounds are tight. These classes have identical strong-fat-shattering and fat-shattering functions.

ACKNOWLEDGMENTS

This research was supported by the Australian Telecommunications and Electronics Research Board and the Australian Research Council. This work was done while Phil Long was visiting the Australian National University and affiliated with Duke University. Phil Long was supported by Air Force Office of Scientific Research Grant F49620-92-J0515. Thanks to Wee Sun Lee, Martin Anthony, and the reviewers for helpful comments.

REFERENCES

1. N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler, Scale-sensitive dimensions, uniform convergence, and learnability, in "Symposium on Foundations of Computer Science, 1993."
2. D. Angluin and L. G. Valiant, Fast probabilistic algorithms for Hamiltonian circuits and matchings, *J. Comput. System Sci.* **18** (1979), 155–193.
3. M. Anthony and P. L. Bartlett, Function learning from interpolation, in "Computational Learning Theory: EUROCOLT '95, 1995."
4. M. Anthony, P. L. Bartlett, Y. Ishai, and J. Shawe-Taylor, Valid generalisation from approximate interpolation, *Combinatorics, Probability and Computing*, to appear.
5. M. Anthony and J. Shawe-Taylor, Valid generalization from approximate interpolation, in "Computational Learning Theory: EUROCOLT '93, 1993."
6. P. Auer, P. M. Long, W. Maass, and G. J. Woeginger, On the complexity of function learning, in "Proceedings, Sixth Annual ACM Conference on Computational Learning Theory, 1993."
7. P. Auer and P. M. Long, Simulating access to hidden information while learning, in "Proceedings, 26th Annual ACM Symposium on the Theory of Computation, 1994."
8. P. L. Bartlett, Learning with a slowly changing distribution, in "Proceedings, Fifth Annual ACM Workshop on Computational Learning Theory, New York, 1992."
9. P. L. Bartlett and P. M. Long, More theorems about scale-sensitive dimensions and learning, in "Proceedings, Eighth Annual ACM Conference on Computational Learning Theory, 1995."
10. P. L. Bartlett, P. M. Long, and R. C. Williamson, Fat-shattering and the learnability of real-valued functions (extended abstract), in "Proceedings, Seventh Annual ACM Conference on Computational Learning Theory, 1994."
11. S. Ben-David, N. Cesa-Bianchi, D. Haussler, and P. Long, Characterizations of learnability for classes of $\{0, \dots, n\}$ -valued functions, *J. Comput. System Sci.* **50** (1995), 74–86.
12. A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth, Learnability and the Vapnik–Chervonenkis dimension, *J. Assoc. Comput. Mach.* **36** (1989), 929–965.
13. S. van de Geer, Regression analysis and empirical processes, in "Centrum voor Wiskunde en Informatica, Amsterdam, 1988."
14. L. Gurvits and P. Koiran, Approximation and learning of convex superpositions, in "Computational Learning Theory: EUROCOLT '95, 1995."
15. D. Haussler, Decision theoretic generalizations of the PAC model for neural net and other learning applications, *Inform. and Comput.* **100** (1992), 78–150.
16. M. J. Kearns and R. E. Schapire, Efficient distribution-free learning of probabilistic concepts (extended abstract), in "Proceedings, 31st Annual Symposium on the Foundations of Computer Science, 1990."
17. M. J. Kearns, R. E. Schapire, and L. M. Sellie, Toward efficient agnostic learning, *Mach. Learning* **17** (1994), 115.
18. N. Merhav and M. Feder, Universal schemes for sequential decision from individual data sequences, *IEEE Trans. Inform. Theory* **39** (1993), 1280–1292.
19. B. K. Natarajan, On learning sets and functions, *Mach. Learning* **4** (1989), 67–97.
20. B. K. Natarajan, Occam's razor for functions, in "Proceedings, Sixth Annual ACM Conference on Computational Learning Theory, 1993."
21. J. K. Patel and C. B. Read, "The Big Book of Facts about the Normal Distribution," Dekker, New York, 1982.
22. D. Pollard, "Convergence of Stochastic Processes," Springer-Verlag, New York, 1984.
23. H. L. Royden, "Real Analysis," Macmillan Co., New York, 1988.
24. J. Shawe-Taylor, M. Anthony, and N. Biggs, Bounding sample size with the Vapnik–Chervonenkis dimension, *Discrete Appl. Math.* **42** (1993), 65–73.
25. H. U. Simon, Bounds on the number of examples needed for learning functions, FB Informatik, LS II, Forschungsbericht Nr. 501, Universität Dortmund, 1993.
26. L. G. Valiant, A theory of the learnable, *Comm. ACM* **27** (1984), 1134–1143.
27. V. N. Vapnik and A. Y. Chervonenkis, On the uniform convergence of relative frequencies of events to their probabilities, *Theory Probab. Appl.* **16** (1971), 264–280.