

# Function Learning from Interpolation

---

MARTIN ANTHONY<sup>1</sup> and PETER L. BARTLETT<sup>2</sup>

<sup>1</sup> Department of Mathematics,  
The London School of Economics and Political Science,  
Houghton Street, London WC2A 2AE, England  
(e-mail: m.anthony@lse.ac.uk)

<sup>2</sup> Department of Systems Engineering,  
Research School of Information Sciences and Engineering,  
The Australian National University, Canberra, 0200 Australia  
(e-mail: Peter.Bartlett@anu.edu.au)

*Received 6 March 1996; revised 26 September 1999*

In this paper, we study a statistical property of classes of real-valued functions that we call approximation from interpolated examples. We derive a characterization of function classes that have this property, in terms of their ‘fat-shattering function’, a notion that has proved useful in computational learning theory. The property is central to a problem of learning real-valued functions from random examples in which we require satisfactory performance from every algorithm that returns a function which approximately interpolates the training examples.

## 1. Introduction

In the problem of learning a real-valued function from examples, a learner sees a sequence of values of an unknown function at a number of randomly chosen points. On the basis of these examples, the learner chooses a function – called a hypothesis – from some class  $\mathcal{H}$  of hypotheses, with the aim that the learner’s hypothesis is close to the target function on future random examples. In this paper we require that, for most training samples, with high probability the absolute difference between the values of the learner’s hypothesis and the target function on a random point is small.

A natural learning algorithm to consider is one that chooses a function in  $\mathcal{H}$  that is close to the target function on the training examples. This poses the following statistical problem: For what function classes  $\mathcal{H}$  will any function in  $\mathcal{H}$  that approximately interpolates the target function on the training examples probably have small absolute

error? More precisely, we have the following definition of approximation from interpolated examples.

**Definition.** Let  $\mathcal{C}, \mathcal{H}$  be sets of functions that map from a set  $X$  to  $\mathbb{R}$ . We say that  $\mathcal{H}$  approximates  $\mathcal{C}$  from interpolated examples if, for all  $\eta, \gamma, \epsilon, \delta \in (0, 1)$ , there is an  $m_0(\eta, \gamma, \epsilon, \delta)$  such that, for every  $t \in \mathcal{C}$  and for every probability measure<sup>†</sup>  $P$  on  $X$ , if  $m \geq m_0(\eta, \gamma, \epsilon, \delta)$ , then with  $P^m$ -probability at least  $1 - \delta$ ,  $\mathbf{x} = (x_1, x_2, \dots, x_m) \in X^m$  has the property that, if  $h \in \mathcal{H}$  and  $|h(x_i) - t(x_i)| < \eta$  for  $1 \leq i \leq m$ , then

$$P(\{x \in X : |h(x) - t(x)| \geq \eta + \gamma\}) < \epsilon.$$

We say that  $m_0(\eta, \gamma, \epsilon, \delta)$  is a *sufficient sample length function* for  $\mathcal{H}$  to approximate  $\mathcal{C}$  from interpolated examples.

Two cases of particular interest are those in which  $\mathcal{C} = \mathcal{H}$  and  $\mathcal{C} = \mathbb{R}^X$ , the set of all functions from  $X$  to  $\mathbb{R}$ . If  $\mathcal{H}$  approximates  $\mathcal{H}$  from interpolated examples, we simply say that  $\mathcal{H}$  approximates from interpolated examples. The main aim of this paper is to find characterizations of classes which approximate from interpolated examples and which approximate  $\mathbb{R}^X$  from interpolated examples.

This problem can be interpreted as a learning problem in which we require satisfactory performance from every algorithm that returns a function that approximately interpolates the training examples. If, instead of requiring that all algorithms in this class be suitable, we require only the existence of a suitable algorithm, no necessary and sufficient conditions on the function class  $\mathcal{H}$  are known. Because an arbitrary amount of information can be conveyed in a single real value, it is possible to construct complicated function classes in which the identity of a function is encoded in its value at every point, and an algorithm can take advantage of this (see [3]). We can avoid this unnatural ‘conspiracy’ between algorithms and function classes in two ways: by requiring that the algorithm be robust in the presence of random observation noise, as was considered in [3], or, contrastingly, by requiring satisfactory performance of every algorithm in a class of reasonable algorithms, as we consider here. Another reason for studying the problem of this paper is that it has implications for learning in the presence of *malicious noise*, in which the labels on the training sample can be any real numbers within  $\eta$  of the true value of the target. This will be discussed later in the paper, but for the moment simply observe that, if  $h$  is  $\beta$ -close to a training sample where the labels have been corrupted to a level of at most  $\beta$ , then  $h$  is certainly  $2\beta$ -close to the target on the sample. If  $\mathcal{H}$  approximates from interpolated examples, we can then deduce that if the sample is large enough then (with high probability)  $h$  is within  $2\beta + \gamma$  of the target on ‘most’ of  $X$ .

Alon, Ben-David, Cesa-Bianchi and Haussler [1] have analysed a model of learning in which the error of a hypothesis is taken to be the expected value of  $(h(x) - t(x))^2$ . Their

<sup>†</sup> More formally, one has a fixed  $\sigma$ -algebra on  $X$ : when  $X$  is countable this is  $2^X$ , and when  $X \subseteq \mathbb{R}^n$ , it is the Borel  $\sigma$ -algebra. Then, by ‘any probability measure on  $X$ ’, we mean ‘for any probability measure on  $\Sigma$ ’, where  $\Sigma$ , the fixed  $\sigma$ -algebra, is understood. The class  $\mathcal{H}$  must have some fairly benign measurability properties; we refer to [12, 8] for details.

results can be used to provide guarantees of small *expected* absolute error. However, the results of this paper provide conditions under which we can (with high probability) have small ‘pointwise’ absolute error *almost everywhere* on the domain, and these results do not follow from those of Alon and co-workers.

In the next section, we define a measure of the complexity of a class  $\mathcal{H}$  of functions (the fat-shattering function), and we state the main result: that the fat-shattering function is the key quantity in this problem. In Sections 3 and 4 we give upper and lower bounds on the number of examples necessary for approximation from interpolated examples. Section 5 describes the implications for learning with malicious noise.

### 2. Definitions and the main result

A number of ways of measuring the ‘expressive power’ of a class  $\mathcal{H}$  of functions have been proposed. This power is quantified by associating a ‘dimension’ to the class. Sometimes this is simply one number depending on  $\mathcal{H}$ . Sometimes – in what is known as a *scale-sensitive dimension* – it is a function depending on  $\mathcal{H}$ .

An important example of the first type of dimension is the *pseudo-dimension* [8, 12]. We say that a finite subset  $S = \{x_1, x_2, \dots, x_d\}$  of  $X$  is *shattered* if there is an  $\mathbf{r} = (r_1, r_2, \dots, r_d) \in \mathbb{R}^d$  such that, for every  $\mathbf{b} = (b_1, b_2, \dots, b_d) \in \{0, 1\}^d$ , there is a function  $h_{\mathbf{b}} \in \mathcal{H}$  with  $h_{\mathbf{b}}(x_i) > r_i$  if  $b_i = 1$  and  $h_{\mathbf{b}}(x_i) < r_i$  if  $b_i = 0$ . The *pseudo-dimension* of  $\mathcal{H}$ , denoted  $\text{Pdim}(\mathcal{H})$ , is the largest cardinality of a shattered set, or infinity if there is no bound on the cardinalities of the shattered sets.

Perhaps the most important scale-sensitive dimension that has been used to date in the development of the theory of learning real-valued functions is the *fat-shattering function*. This is a scale-sensitive version of the pseudo-dimension and was introduced by Kearns and Schapire [9]. Suppose that  $\mathcal{H}$  is a set of functions from  $X$  to  $[0, 1]$  and that  $\gamma \in (0, 1)$ . We say that a finite subset  $S = \{x_1, x_2, \dots, x_d\}$  of  $X$  is  $\gamma$ -*shattered* if there is an  $\mathbf{r} = (r_1, r_2, \dots, r_d) \in \mathbb{R}^d$  such that, for every  $\mathbf{b} = (b_1, b_2, \dots, b_d) \in \{0, 1\}^d$ , there is a function  $h_{\mathbf{b}} \in \mathcal{H}$  with  $h_{\mathbf{b}}(x_i) \geq r_i + \gamma$  if  $b_i = 1$  and  $h_{\mathbf{b}}(x_i) \leq r_i - \gamma$  if  $b_i = 0$ . Thus,  $S$  is  $\gamma$ -shattered if it is shattered with a ‘width of shattering’ of at least  $\gamma$ . We define the *fat-shattering function*,  $\text{fat}_{\mathcal{H}} : \mathbb{R}^+ \rightarrow \mathbb{N}_0 \cup \{\infty\}$ , as

$$\text{fat}_{\mathcal{H}}(\gamma) = \max \{|S| : S \subseteq X \text{ is } \gamma\text{-shattered by } \mathcal{H}\},$$

or  $\text{fat}_{\mathcal{H}}(\gamma) = \infty$  if the maximum does not exist. (Here,  $\mathbb{N}_0$  denotes the set of nonnegative integers.) It is easy to see that  $\text{Pdim}(\mathcal{H}) = \lim_{\gamma \rightarrow 0} \text{fat}_{\mathcal{H}}(\gamma)$ . It should be noted, however, that it is possible for the pseudo-dimension to be infinite, even when  $\text{fat}_{\mathcal{H}}(\gamma)$  is finite for all  $\gamma$ . We shall say that  $\mathcal{H}$  has *finite fat-shattering function* whenever it is the case that, for all  $\gamma \in (0, 1)$ ,  $\text{fat}_{\mathcal{H}}(\gamma)$  is finite.

The fat-shattering function plays an important role in the learning theory of real-valued functions. Kearns and Schapire [9] proved that if a class of probabilistic concepts is learnable, then the class has finite fat-shattering function. (A probabilistic concept  $f$  is a  $[0, 1]$ -valued function. In this model, the learner sees examples  $(x_i, y_i)$ , where  $\Pr(y_i = 1) = f(x_i)$ .) Alon, Ben-David, Cesa-Bianchi and Haussler [1] proved, conversely, that, if a class of probabilistic concepts has finite fat-shattering function, then it is

learnable. The main result in [3] is that finiteness of the fat-shattering function of a class of  $[0, 1]$ -valued functions is a necessary and sufficient condition for learning with random observation noise.

Our main result is the following.

**Theorem 2.1.** *Suppose that  $\mathcal{H}$  is a set of functions from a set  $X$  to  $[0, 1]$ . Then the following propositions are equivalent.*

- (1)  $\mathcal{H}$  approximates from interpolated examples.
- (2)  $\mathcal{H}$  approximates  $\mathbb{R}^X$  from interpolated examples.
- (3)  $\mathcal{H}$  has finite fat-shattering function.

### 3. The upper bound

In this section, we prove that finite fat-shattering function is a sufficient condition for approximation from interpolated examples and we provide a suitable sample length function  $m_0(\eta, \gamma, \epsilon, \delta)$ .

We first need the notion of *covering numbers*  $\mathcal{N}_A(\alpha, d)$ , as used extensively in [8, 1, 6], for instance. Suppose that  $(A, d)$  is a pseudo-metric space and  $\alpha > 0$ . Then, a subset  $N$  of  $A$  is said to be an  $\alpha$ -cover for a subset  $B$  of  $A$  if, for every  $x \in B$ , there is an  $\hat{x} \in N$  such that  $d(x, \hat{x}) \leq \alpha$ . The metric space is *totally bounded* if there is a finite  $\alpha$ -cover for  $A$ , for all  $\alpha > 0$ . When  $(A, d)$  is totally bounded, we shall denote the minimal cardinality of an  $\alpha$ -cover for  $A$  by  $\mathcal{N}_A(\alpha, d)$  for  $\alpha > 0$ . A subset  $M$  of  $A$  is said to be  $\alpha$ -separated if, for all distinct  $x, y \in M$ ,  $d(x, y) \geq \alpha$ . We shall denote the maximal cardinality of an  $\alpha$ -separated subset of  $A$  by  $\mathcal{M}_A(\alpha, d)$ . It is easy to show that

$$\mathcal{M}_A(2\alpha, d) \leq \mathcal{N}_A(\alpha, d) \leq \mathcal{M}_A(\alpha, d)$$

(see [10]), so  $\mathcal{M}_A(\alpha, d)$  is always defined if  $(A, d)$  is totally bounded. Suppose now that  $\mathcal{H}$  is a set of functions from a set  $X$  to  $[0, 1]$  and that  $\mathbf{x} = (x_1, x_2, \dots, x_m) \in X^m$ , where  $m$  is a positive integer. We may define a pseudo-metric  $l_{\mathbf{x}}^{\infty}$  on  $\mathcal{H}$  as follows: for  $g, h \in \mathcal{H}$ ,

$$l_{\mathbf{x}}^{\infty}(g, h) = \max_{1 \leq i \leq m} |g(x_i) - h(x_i)|.$$

(This metric has been used in [6, 1], for example.) Alon, Ben-David, Cesa-Bianchi and Haussler [1] obtained (essentially) the following result bounding the  $l_{\mathbf{x}}^{\infty}$ -covering number of  $\mathcal{H}$  in terms of the fat-shattering function of  $\mathcal{H}$ .

**Lemma 3.1.** *Suppose that  $\mathcal{H}$  is a set of functions from  $X$  to  $[0, 1]$  and that  $\mathcal{H}$  has finite fat-shattering function. Let  $m \in \mathbb{N}$ , and  $\mathbf{x} \in X^m$ . Then the pseudo-metric space  $(\mathcal{H}, l_{\mathbf{x}}^{\infty})$  is totally bounded. Suppose  $\alpha > 0$ . Let  $d = \text{fat}_{\mathcal{H}}(\alpha/4)$  and*

$$y = \sum_{i=1}^d \binom{m}{i} \left( \left\lceil \frac{2}{\alpha} \right\rceil \right)^i.$$

Then, provided  $m \geq \log y + 1$ ,

$$\mathcal{N}_{\mathcal{H}}(\alpha, l_{\mathbf{x}}^{\infty}) < 2 \left( m \left\lceil \frac{2}{\alpha} \right\rceil^2 \right)^{\log y}.$$

Here, as elsewhere in the paper,  $\log$  denotes logarithm to base 2. We then have the following result.

**Theorem 3.2.** *Suppose that  $\mathcal{H}$  is a class of functions mapping from a domain  $X$  to the real interval  $[0, 1]$  and that  $\mathcal{H}$  has finite fat-shattering function. Let  $t$  be any function from  $X$  to  $\mathbb{R}$  and let  $\gamma, \eta, \epsilon > 0$ . Let  $P$  be any probability distribution on  $X$  and define  $B$  to be the set of functions  $h \in \mathcal{H}$  for which  $P(\{x \in X : |h(x) - t(x)| \geq \eta + \gamma\}) \geq \epsilon$ . Let  $d = \text{fat}_{\mathcal{H}}(\gamma/8)$  and let*

$$y = \sum_{i=1}^d \binom{2m}{i} \left( \left\lceil \frac{4}{\gamma} \right\rceil \right)^i.$$

Then, for  $m \geq \max(8/\epsilon, \log y + 1)$ , the probability that some  $h$  in  $B$  has  $|h(x_i) - t(x_i)| < \eta$  for  $1 \leq i \leq m$  is at most

$$4 \left( 2m \left\lceil \frac{4}{\gamma} \right\rceil^2 \right)^{\log y} 2^{-\epsilon m/2}.$$

**Proof.** The proof is based on a technique analogous to that used in [13, 5, 8], where we ‘symmetrize’ and then ‘combinatorially bound’. The first step – symmetrization – relates the desired probability to a ‘sample-based’ one. Fix  $t, P, m$ , the parameters  $\gamma, \eta, \epsilon$ , and hence the set  $B$ . It is easy to show using standard techniques that

$$P^m \{ \mathbf{x} \in X^m : \exists h \in B, |h(x_i) - t(x_i)| < \eta \ (1 \leq i \leq m) \} \leq P^{2m}(R),$$

where

$$R = \{ \mathbf{xy} \in X^{2m} : \exists h \in B, |h(x_i) - t(x_i)| < \eta \ (1 \leq i \leq m) \\ \text{and } |\{i : |h(y_i) - t(y_i)| \geq \eta + \gamma\}| > \epsilon m/2 \},$$

and  $\mathbf{xy} \in X^{2m}$  denotes the concatenation of  $\mathbf{x}, \mathbf{y} \in X^m$ .

The next step is to bound the probability of  $R$  using combinatorial techniques. For this, let  $\Gamma$  be the ‘swapping group’ [12] of permutations on the set  $\{1, 2, \dots, 2m\}$ . This is the group generated by the transpositions  $(i, m+i)$  for  $1 \leq i \leq m$ . The group  $\Gamma$  acts in a natural way on vectors in  $X^{2m}$ : for  $\sigma \in \Gamma$  and  $\mathbf{z} \in X^{2m}$ , we define  $\sigma\mathbf{z}$  to be

$$(z_{\sigma(1)}, z_{\sigma(2)}, \dots, z_{\sigma(2m)}).$$

Let  $\Gamma(R, \mathbf{z}) = |\{\sigma \in \Gamma : \sigma\mathbf{z} \in R\}|$  be the number of permutations in  $\Gamma$  taking  $\mathbf{z}$  into  $R$ . It is well known that, since  $P^{2m}$  is a product distribution, we have

$$P^{2m}(R) \leq \frac{1}{2^m} \max_{\mathbf{z} \in X^{2m}} \Gamma(R, \mathbf{z}).$$

Now, let us fix  $\mathbf{z} \in X^{2m}$  and consider the pseudo-metric space  $(B, l_{\mathbf{z}}^{\infty})$ . Since  $\mathcal{H}$  has finite

fat-shattering function, so does  $B$  and Lemma 3.1 implies that this pseudo-metric space is totally bounded. Let  $N = \{\hat{h}_1, \hat{h}_2, \dots, \hat{h}_n\}$  be a minimal  $\gamma/2$ -cover for  $B$ . From Lemma 3.1,

$$n < 2 \left( 2m \left\lceil \frac{4}{\gamma} \right\rceil^2 \right)^{\log y}.$$

Since  $N$  is a  $\gamma/2$ -cover, given  $h \in B$ , there is an  $\hat{h} \in N$  such that  $l_z^\infty(\hat{h}, h) < \gamma/2$ , which means that, for  $1 \leq i \leq 2m$ ,  $|\hat{h}(z_i) - h(z_i)| < \gamma/2$ . Suppose that  $\sigma z = \mathbf{xy} \in R$ . Then, by the definition of  $R$ , there is some  $h \in B$  such that  $|h(x_i) - t(x_i)| < \eta$  for  $1 \leq i \leq m$  and such that, for more than  $\epsilon m/2$  of the  $y_i$ ,  $|h(y_i) - t(y_i)| \geq \eta + \gamma$ . But (taking  $\hat{h}$  to be, as described above, a function in the cover  $\gamma/2$ -close to  $h$ ) this implies that there is an  $\hat{h} \in N$  such that, for  $1 \leq i \leq m$ ,  $|\hat{h}(x_i) - t(x_i)| < \eta + \gamma/2$ , and such that, for more than  $\epsilon m/2$  of the  $y_i$ ,  $|\hat{h}(y_i) - t(y_i)| \geq \eta + \gamma/2$ . It follows from this that, if  $\sigma z \in R$ , then, for some  $l$  between 1 and  $n$ ,  $\sigma z$  belongs to the set  $\hat{R}_l$ , defined by

$$\begin{aligned} \hat{R}_l = \{ \mathbf{xy} \in X^{2m} : & |\hat{h}_l(x_i) - t(x_i)| < \eta + \gamma/2 (1 \leq i \leq m), \\ & \text{and } |\{i : |\hat{h}_l(y_i) - t(y_i)| \geq \eta + \gamma/2\}| > \epsilon m/2 \}. \end{aligned}$$

Let  $\Gamma(\hat{R}_l, \mathbf{z})$  be the number of  $\sigma$  in  $\Gamma$  for which  $\sigma z \in \hat{R}_l$ . Since  $\sigma z \in R$  implies  $\sigma z \in \hat{R}_l$  for some  $l$ , we have

$$\Gamma(R, \mathbf{z}) \leq \sum_{l=1}^n \Gamma(\hat{R}_l, \mathbf{z}).$$

Consider a particular  $l$  between 1 and  $n$  and suppose that  $\Gamma(\hat{R}_l, \mathbf{z}) \neq 0$ . Let  $k$  be the number of indices  $i$  between 1 and  $2m$  such that  $|\hat{h}_l(z_i) - t(z_i)| \geq \eta + \gamma/2$ . Then  $\epsilon m/2 < k \leq m$ . The number of permutations  $\sigma$  in  $\Gamma$  for which  $\sigma z$  belongs to  $\hat{R}_l$  is then equal to  $2^{m-k}$ , which is less than  $2^{m(1-\epsilon/2)}$ . (The  $z_i$  which can be ‘swapped’ are precisely those  $m - k$  satisfying  $|\hat{h}_l(z_{m+i}) - t(z_{m+i})| < \eta + \gamma/2$ .) It follows that

$$P^{2m}(R) < \frac{1}{2^m} \sum_{i=1}^n 2^{m(1-\epsilon/2)} \leq n 2^{-\epsilon m/2} \leq 2 \left( 2m \left\lceil \frac{4}{\gamma} \right\rceil^2 \right)^{\log y} 2^{-\epsilon m/2}.$$

The statement of the theorem now follows. □

We thus obtain the following corollary, which shows that finiteness of  $\text{fat}_{\mathcal{H}}$  implies that  $\mathcal{H}$  approximates  $\mathbb{R}^X$  from interpolated examples, and hence  $\mathcal{H}$  approximates  $\mathcal{H}$  from interpolated examples. The proof is an easy calculation.

**Corollary 3.3.** *Suppose that  $\mathcal{H}$  is a set of functions from  $X$  to  $[0, 1]$  and that  $\mathcal{H}$  has finite fat-shattering function. Then  $\mathcal{H}$  approximates  $\mathbb{R}^X$  from interpolated examples. Furthermore, there is a positive constant  $K$  such that a sufficient sample length function is*

$$m_0(\gamma, \eta, \epsilon, \delta) = \frac{K}{\epsilon} \left( \log \left( \frac{1}{\delta} \right) + d \log^2 \left( \frac{d}{\gamma \epsilon} \right) \right),$$

where  $d = \text{fat}_{\mathcal{H}}(\gamma/8)$ . □

4. The lower bound

In this section, we give lower bounds on the number of examples necessary for  $\mathcal{H}$  to approximate  $\mathbb{R}^X$  from interpolated examples and for  $\mathcal{H}$  to approximate  $\mathcal{H}$  from interpolated examples. The bounds are in terms of  $\text{fat}_{\mathcal{H}}$ , the fat-shattering function of  $\mathcal{H}$ . To prove them, we consider a discretized version of  $\mathcal{H}$ . We then consider a number of notions of dimension for classes that map to these discrete sets, and show that a large family of these dimensions consists of closely related members. This family includes a version of the Natarajan dimension – see [11] – for which it is easy to prove lower bounds. Since the fat-shattering function is also a member of this family of closely related dimension, we obtain the lower bound. This broad outline is similar to the approach adopted by Ben-David, Cesa-Bianchi, Haussler and Long [4], who consider learning  $[n]$ -valued functions.

We first define the discretization we shall use. For  $a \in [0, 1]$ , let  $D_\gamma(a) = \lceil a/\gamma \rceil$ . For a function  $f : X \rightarrow [0, 1]$ , let  $D_\gamma(f) : X \rightarrow \{0, 1, \dots, \lceil 1/\gamma \rceil\}$  be defined as the composition of  $D_\gamma$  and  $f$ . Let  $D_\gamma(\mathcal{H})$  denote  $\{D_\gamma(f) : f \in \mathcal{H}\}$ . Functions in  $D_\gamma(\mathcal{H})$  map to  $\{0, 1, \dots, n\}$ , where  $n = \lceil 1/\gamma \rceil$ . Let  $[n]$  denote  $\{0, 1, \dots, n\}$ .

From the definition of the fat-shattering function,

$$\text{fat}_{\mathcal{H}}(\alpha) \leq \text{fat}_{D_\gamma(\mathcal{H})} \left( \frac{1}{2} \left\lfloor \frac{2\alpha}{\gamma} \right\rfloor \right)$$

for  $\alpha, \gamma \in \mathbb{R}^+$ .

We consider the following notions of dimension, defined using classes of  $\{0, 1, *\}$ -valued functions on  $[n]$ .

**Definition.** If  $\mathcal{F}$  is a class of  $[n]$ -valued functions defined on  $X$  and  $\Psi$  is a class of  $\{0, 1, *\}$ -valued functions defined on  $[n]$ , we say that  $\mathcal{F}$   $\Psi$ -shatters  $\mathbf{x} = (x_1, \dots, x_d) \in X^d$  if there is a sequence  $(\phi_1, \dots, \phi_d) \in \Psi^d$  such that

$$\{0, 1\}^d \subseteq \{(\phi_1(f(x_1)), \dots, \phi_d(f(x_d))) : f \in \mathcal{F}\}.$$

The  $\Psi$ -dimension of  $\mathcal{F}$ , denoted  $\Psi\text{-dim}(\mathcal{F})$ , is the size of the largest  $\Psi$ -shattered sequence, or infinity if there is no largest sequence.

We can express the fat-shattering function  $\text{fat}_{\mathcal{F}}(k)$  as a dimension of this type, for  $k \geq 2$ . Define  $\Psi_{\text{fat}}(k) = \{\psi_i : i \in \{0, \dots, n - k\}\}$ , with

$$\psi_i(z) = \begin{cases} 1, & z \geq i + k, \\ *, & i < z < i + k, \\ 0, & z \leq i, \end{cases}$$

for  $z \in [n]$ . Then  $\text{fat}_{\mathcal{F}}(\alpha) = \Psi_{\text{fat}}(\lceil 2\alpha \rceil)\text{-dim}(\mathcal{F})$  for all classes  $\mathcal{F}$  of functions from  $X$  to  $[n]$ .

The  $\Psi_{\text{gnat}}$ -dimension (a ‘gapped’ version of the Natarajan dimension) will be useful, since it is easy to prove lower bounds using this dimension.

**Definition.** Let  $\Psi_{\text{gnat}}(k)$  be the following set of  $\{0, 1, *\}$ -valued functions defined on  $[n]$ , where  $k \in \{2, 3, \dots\}$ :

$$\Psi_{\text{gnat}}(k) = \{\psi_{i,j} : i, j \in [n], |i - j| \geq k\},$$

with

$$\psi_{i,j}(\alpha) = \begin{cases} 1, & \alpha = i, \\ 0, & \alpha = j, \\ *, & \text{otherwise.} \end{cases}$$

The following result and its proof (which we omit) are similar to the key result in [1], Lemma 15, which bounds covering numbers of  $\mathcal{F}$  in terms of  $\text{fat}_{\mathcal{F}}(1)$ . We will see later that it generalizes that lemma, since the  $\Psi_{\text{gnat}}(2)$ -dimension is the smallest of a family of dimensions that includes  $\text{fat}_{\mathcal{F}}(1)$ .

**Lemma 4.1.** Let  $k \geq 2$  and  $n \geq 1$  be integers. Suppose that  $\mathcal{F}$  is a class of  $[n]$ -valued functions defined on  $X$  satisfying  $\Psi_{\text{gnat}}(k)\text{-dim}(\mathcal{F}) \leq d$ , and that  $m \geq \log y + 1$ , where

$$y = \sum_{i=0}^d \binom{m}{i} n^{2i}.$$

Then

$$\max_{\mathbf{x} \in X^m} \mathcal{M}_{\mathcal{F}}(\mathbf{x}, l_{\mathbf{x}}^{\infty}) < 2(2mn^2)^{\log y}.$$

The ‘ $k$ -gapped distinguishers’ correspond to a family of dimensions that includes the  $\Psi_{\text{gnat}}(k)$ -dimension and the fat-shattering function at a certain scale.

**Definition.** Let  $k \geq 2$ . A set  $\Psi$  of functions from  $[n]$  to  $\{0, 1, *\}$  is a  $k$ -gapped distinguisher if it satisfies:

- (1) for all  $i \in \{0, 1, \dots, n - k\}$  and  $j \in \{i + k, \dots, n\}$ , there is a function  $\psi \in \Psi$  and a bit  $b \in \{0, 1\}$  such that  $\psi(i) = b$  and  $\psi(j) = 1 - b$ ;
- (2)  $\min \{|i - j| : i, j \in [n], \exists \psi \in \Psi, \psi(i) = 0, \psi(j) = 1\} = k$ .

In addition to the set  $\Psi_{\text{gnat}}(k)$ , another important example of a  $k$ -gapped distinguisher is the class

$$\Psi_g(k) = \{\psi \in \{0, 1, *\}^{[n]} : \min \{|i - j| : i, j \in [n], \psi(i) = 0, \psi(j) = 1\} = k\}.$$

In fact  $\Psi_g(k)$  is the largest  $k$ -gapped distinguisher, in the sense that it contains any other  $k$ -gapped distinguisher.

**Lemma 4.2.** Suppose  $\mathcal{F}$  is a class of  $[n]$ -valued functions defined on  $X$ ,  $\Psi$  is a class of  $\{0, 1, *\}$ -valued functions defined on  $[n]$ , and  $k \geq 2$ . If  $\Psi$  is a  $k$ -gapped distinguisher then  $\Psi_{\text{gnat}}(k)\text{-dim}(\mathcal{F}) \leq \Psi\text{-dim}(\mathcal{F}) \leq \Psi_g(k)\text{-dim}(\mathcal{F})$ .



**Proof.** Take a  $\Psi_{\text{gnat}}(k)$ -shattered sequence  $\mathbf{x} \in X^d$ . Since  $\Psi$  is a  $k$ -gapped distinguisher, for all  $\psi_{i,j} \in \Psi_{\text{gnat}}(k)$  there is a  $\phi \in \Psi$  and a  $b \in \{0, 1\}$  for which  $\phi(i) = b$  and  $\phi(j) = 1 - b$ . It follows that  $\mathbf{x}$  is  $\Psi$ -shattered, which gives the first inequality. The second inequality follows from the fact that  $\Psi \subseteq \Psi_g(k)$ .  $\square$

It follows from Lemma 4.2 that Lemma 4.1 generalizes Alon, Ben-David, Cesa-Bianchi and Haussler’s Lemma 15 [1], which gave a similar result for the  $\Psi_{\text{fat}}(2)$ -dimension. The following result shows that the  $\Psi_{\text{gnat}}(k)$ -dimension, the  $\Psi_g(k)$ -dimension, and the  $\Psi$ -dimension (for any  $k$ -gapped distinguisher  $\Psi$ ) are all closely related.

**Lemma 4.3.** *Let  $k \geq 2$ . Let  $\mathcal{F}$  be a class of functions that map from  $X$  to  $[n]$ , satisfying  $\Psi_g(k)\text{-dim}(\mathcal{F}) \geq d \geq 2$ . Then*

$$\Psi_{\text{gnat}}(k)\text{-dim}(\mathcal{F}) > \frac{d}{3 \log^2(2dn^2)}.$$

**Proof.** Suppose  $\mathbf{x} = (x_1, \dots, x_d) \in X^d$  is  $\Psi_g(k)$ -shattered by  $\mathcal{F}$ . The definition of  $\Psi_g(k)$  implies that any (minimal) subset of  $\mathcal{F}$  that  $\Psi_g(k)$ -shatters  $\mathbf{x}$  is  $k$ -separated, so  $\mathcal{M}_{\mathcal{F}}(k, l_{\mathbf{x}}^{\infty}) \geq 2^d$ .

Suppose  $\Psi_g(k)\text{-dim}(\mathcal{F}) = d_N$ , and let

$$y = \sum_{i=0}^{d_N} \binom{d}{i} n^{2i}.$$

If  $d > \log y$  then, by Lemma 4.1,

$$2^d \leq \mathcal{M}_{\mathcal{F}}(k, l_{\mathbf{x}}^{\infty}) < 2(2dn^2)^{\log y},$$

so

$$d < 1 + \log y \log(2dn^2). \tag{4.1}$$

Alternatively, if  $d \leq \log y$ , then (4.1) is obviously true. Clearly,  $d_N = 0$  only if  $d = 0$ , so assume  $d_n \geq 1$ . Then  $y \leq 2d_N d^{d_N} n^{2d_N}$ , so we have

$$\begin{aligned} d &< 1 + \log y \log(2dn^2) \\ &\leq 1 + (\log(2d_N) + d_N \log(dn^2)) \log(2dn^2) \\ &\leq 3d_N \log^2(2dn^2). \end{aligned} \quad \square$$

The following result follows easily from [7, Theorem 1], which gives a lower bound on the number of examples necessary for learning  $\{0, 1\}^{[d]}$  in the probably approximately correct model (see also [5]).

**Lemma 4.4.** *Let  $0 < \epsilon \leq 1/8$ ,  $0 < \delta < 1/100$ , and  $d \geq 1$ . If*

$$m < \max \left( \frac{d}{32\epsilon}, \frac{1 - \epsilon}{\epsilon} \ln \frac{1}{\delta} \right),$$

then there is a distribution  $P$  on  $[d]$  and a function  $t \in \{0, 1\}^{[d]}$  such that

$$P^m \{ \mathbf{x} \in X^m : \exists f \in \{0, 1\}^{[d]} \text{ such that } f(x_i) = t(x_i), i = 1, \dots, m \text{ and } P \{ y : f(y) \neq t(y) \} \geq \epsilon \} \geq \delta.$$

We use Lemma 4.4 to prove the following lower bound on the sample length function for approximating  $\mathbb{R}^X$  from interpolated examples.

**Theorem 4.5.** *Suppose  $\mathcal{H}$  is a class of  $[0, 1]$ -valued functions defined on a set  $X$ ,  $0 < \gamma < \eta < 1$ , and  $\epsilon, \delta \in (0, 1)$ . Then if  $\text{fat}_{\mathcal{H}}(\gamma) \geq d \geq 1$  and  $\gamma^2 \geq 4d2^{-\sqrt{d}/6}$ , any sample length function  $m_0$  for  $\mathcal{H}$  to approximate  $\mathbb{R}^X$  from interpolated examples satisfies*

$$m_0(\eta, \gamma, \epsilon, \delta) \geq \max \left( \frac{1}{32\epsilon} \left( \frac{d}{3 \log^2(4d/\gamma^2)} - 1 \right), \frac{1}{\epsilon} \log \frac{1}{\delta} \right).$$

**Proof.** Fix  $0 < \gamma < \eta < 1$ , define  $n = \lceil 1/\gamma \rceil$ , and suppose  $\text{fat}_{\mathcal{H}}(\gamma) \geq d$ . Let  $\mathcal{F} = D_\gamma(\mathcal{H})$ . Then  $\text{fat}_{\mathcal{F}}(1) \geq d$ , so  $\Psi_{\text{gnat}}(2)\text{-dim}(\mathcal{F}) > k$ , where  $k = d/(3 \log^2(2dn^2))$ . Consider a sequence  $(x_1, \dots, x_k) \in X^k$  that is  $\Psi_{\text{gnat}}(2)$ -shattered by  $\mathcal{F}$ . Clearly, there is a subset  $\mathcal{H}_0 \subseteq \mathcal{H}$  with  $|\mathcal{H}_0| = 2^k$  and a sequence  $(\psi_{a_1, b_1}, \dots, \psi_{a_k, b_k}) \in \Psi_{\text{gnat}}(2)^k$  such that

$$\{ (\psi_{a_1, b_1}(f(x_1)), \dots, \psi_{a_k, b_k}(f(x_k))) : f \in D_\gamma(\mathcal{H}_0) \} = \{0, 1\}^k.$$

Without loss, we can assume that  $a_j > b_j$  for  $j = 1, \dots, k$ .

Now, if  $m < \max((k-1)/(32\epsilon), ((1-\epsilon)/\epsilon) \ln(1/\delta))$  and  $k \geq 2$  (for which  $\gamma^2 \geq 4d2^{-\sqrt{d}/6}$  suffices), Lemma 4.4 implies that there is a distribution  $P$  on  $\{1, \dots, k\}$  and a function  $p : \{1, \dots, k\} \rightarrow \{0, 1\}$  such that

$$P^m \{ l \in \{1, \dots, k\}^m : \exists p' : \{1, \dots, k\} \rightarrow \{0, 1\} \text{ such that } p(l_i) = p'(l_i) \text{ for } i = 1, \dots, m \text{ and } P \{ y \in \{1, \dots, k\} : p(y) \neq p'(y) \} \geq \epsilon \} \geq \delta.$$

Choose a function  $t : X \rightarrow \mathbb{R}$  satisfying

$$t(x_j) = \begin{cases} (a_j - 1)\gamma + \eta, & p(j) = 1, \\ b_j\gamma - \eta + \Delta, & p(j) = 0, \end{cases}$$

for  $j = 1, \dots, k$ , where

$$\Delta = \frac{1}{2} \min \{ h(x_j) - (a_j - 1)\gamma : h \in \mathcal{H}_0 \text{ and } h(x_j) > (a_j - 1)\gamma, j = 1, \dots, k \}.$$

For each function  $h \in \mathcal{H}_0$  define  $f_h = D_\gamma(h)$ . Let  $p_h : \{1, \dots, k\} \rightarrow \{0, 1\}$  be defined by

$$p_h(j) = \begin{cases} 1, & f_h(x_j) = a_j, \\ 0, & f_h(x_j) = b_j. \end{cases}$$

Clearly, if  $|h(x_j) - t(x_j)| < \eta$  for some  $h \in \mathcal{H}_0$  and some  $j \in \{1, \dots, k\}$ , then  $h(x_j) \in ((a_j - 1)\gamma, a_j\gamma] \cup ((b_j - 1)\gamma, b_j\gamma]$  so  $p_h(j) = p(j)$ . Also, if  $p_h(j) \neq p(j)$  for some  $h \in \mathcal{H}$ , then  $|h(x_j) - t(x_j)| > \eta + \gamma$ . It follows that  $P \{ y \in \{1, \dots, k\} : p_h(y) \neq p(y) \} \geq \epsilon$  implies  $Q \{ y \in X : |h(y) - t(y)| \geq \eta + \gamma \} \geq \epsilon$ , where  $Q$  is the discrete probability distribution on  $X$

satisfying  $Q(x_j) = P(j)$  for  $j = 1, \dots, k$ . So

$$\begin{aligned}
 & Q^m \{ \mathbf{y} \in X^m : \exists h \in \mathcal{H} \text{ such that } |h(y_i) - t(y_i)| < \eta \text{ for } i = 1, \dots, m \text{ and} \\
 & \quad Q \{ \mathbf{y} \in X : |h(y) - t(y)| \geq \eta + \gamma \} \geq \epsilon \} \geq \\
 & P^m \{ \mathbf{y} \in X^m : \exists h \in \mathcal{H}_0 \text{ such that } |h(y_i) - t(y_i)| < \eta \text{ for } i = 1, \dots, m \text{ and} \\
 & \quad Q \{ \mathbf{y} \in X : |h(y) - t(y)| \geq \eta + \gamma \} \geq \epsilon \} \geq \\
 & P^m \{ l \in \{1, \dots, k\}^m : \exists p' : \{1, \dots, k\} \rightarrow \{0, 1\} \text{ such that} \\
 & \quad p(l_i) = p'(l_i) \text{ for } i = 1, \dots, m \text{ and} \\
 & \quad P \{ y \in \{1, \dots, k\} : p(y) \neq p'(y) \} \geq \epsilon \} \geq \delta. \quad \square
 \end{aligned}$$

We also have the following result which bounds from below the sample length function for  $\mathcal{H}$  to approximate  $\mathcal{H}$  from interpolated examples.

**Theorem 4.6.** *Suppose  $\mathcal{H}$  is a class of  $[0, 1]$ -valued functions defined on a set  $X$ ,  $0 < \gamma < 1$ ,  $3\gamma/2 \leq \eta < 1$ , and  $\epsilon, \delta \in (0, 1)$ . If  $d$  satisfies  $\text{fat}_{\mathcal{H}}(\eta + \gamma) \geq d \geq 1$  and  $\gamma^2 \geq 4d2^{-\sqrt{d/6}}$ , then any sample length function  $m_0$  for  $\mathcal{H}$  to approximate  $\mathcal{H}$  from interpolated examples satisfies*

$$m_0(\eta, \gamma, \epsilon, \delta) \geq \max \left( \frac{1}{32\epsilon} \left( \frac{d}{3 \log^2(4d/\gamma^2)} - 1 \right), \frac{1}{\epsilon} \log \frac{1}{\delta} \right).$$

**Proof.** Fix  $0 < \gamma < 1$  and  $3\gamma/2 \leq \eta < 1$ , define  $n = \lceil 1/\gamma \rceil$ , and suppose  $d \leq \text{fat}_{\mathcal{H}}(\eta + \gamma)$ . Let  $\mathcal{F} = D_\gamma(\mathcal{H})$ . Then

$$\text{fat}_{\mathcal{F}} \left( \frac{1}{2} \left\lfloor \frac{2(\eta + \gamma)}{\gamma} \right\rfloor \right) \geq d,$$

so  $\Psi_{\text{fat}}(\lfloor 2\eta/\gamma \rfloor + 1)\text{-dim}(\mathcal{F}) \geq d$ , hence  $\Psi_{\text{gnat}}(\lfloor 2\eta/\gamma \rfloor + 1)\text{-dim}(\mathcal{F}) > k$ , where  $k = d/(3 \log^2(2dn^2))$ . Consider a sequence  $(x_1, \dots, x_k) \in X^k$  that is  $\Psi_{\text{gnat}}(\lfloor 2\eta/\gamma \rfloor + 1)$ -shattered by  $\mathcal{F}$ . Clearly, there is a subset  $\mathcal{H}_0 \subseteq \mathcal{H}$  with  $|\mathcal{H}_0| = 2^k$  and a sequence  $(\psi_{a_1, b_1}, \dots, \psi_{a_k, b_k}) \in \Psi_{\text{gnat}}(\lfloor 2\eta/\gamma \rfloor + 1)^k$  such that

$$\{ (\psi_{a_1, b_1}(f(x_1)), \dots, \psi_{a_k, b_k}(f(x_k))) : f \in D_\gamma(\mathcal{H}_0) \} = \{0, 1\}^k.$$

Fix a function  $t \in \mathcal{H}_0$ . Any function  $h \in \mathcal{H}$  that has

$$\psi_{a_i, b_i}(D_\gamma(h)(x_i)) = \psi_{a_i, b_i}(D_\gamma(t)(x_i))$$

satisfies  $|h(x_i) - t(x_i)| < \gamma < \eta$ . Any function  $h$  in  $\mathcal{H}$  that has

$$\psi_{a_i, b_i}(D_\gamma(h)(x_i)) \neq \psi_{a_i, b_i}(D_\gamma(t)(x_i))$$

satisfies

$$\begin{aligned}
 |h(x_i) - t(x_i)| & \geq \left\lfloor \frac{2\eta}{\gamma} \right\rfloor \gamma \\
 & = 2\gamma + \left\lfloor \frac{2(\eta - \gamma)}{\gamma} \right\rfloor \gamma \\
 & \geq 2\gamma + \eta - \gamma = \eta + \gamma,
 \end{aligned}$$

since  $(\eta - \gamma)/\gamma \geq 1/2$  and  $\lfloor 2\alpha \rfloor \geq \alpha$  for  $\alpha \geq 1/2$ .

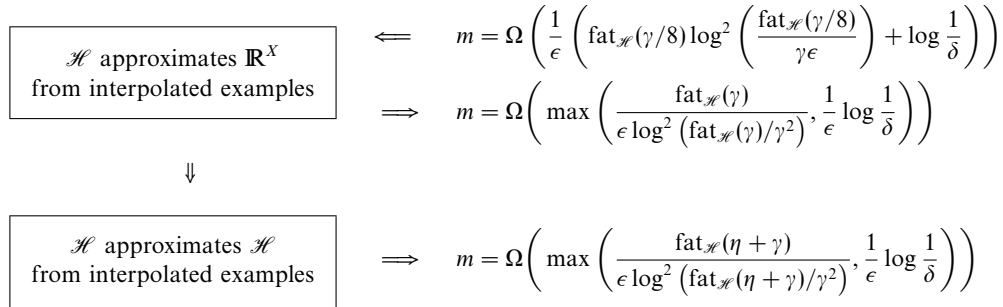


Figure 1 Sample complexity bounds

Using the same argument as in the proof of Theorem 4.5, there is a distribution  $P$  on  $X$  such that if  $m$  is too small then, with  $P^m$ -probability at least  $\delta$ , some  $h \in \mathcal{H}$  is within  $\eta$  of  $t$  on a random sample, but  $P(|h - t| \geq \eta + \gamma) \geq \epsilon$ .  $\square$

5. Discussion

Figure 1 shows the sample complexity bounds for approximation from interpolated examples. (The diagram omits the requirement in the lower bounds that  $\gamma$  is not too small as a function of  $\text{fat}_{\mathcal{H}}(\gamma)$  and  $\text{fat}_{\mathcal{H}}(\eta + \gamma)$ .) These bounds imply Theorem 2.1.

Notice that the upper and lower bounds on the sample length for  $\mathcal{H}$  to approximate  $\mathbb{R}^X$  from interpolated examples are within log factors of each other. These sample complexity bounds are also relevant to the problem of learning real-valued functions in the presence of malicious noise. Suppose a learner sees a sequence of training examples that correspond to the values of a target function corrupted with arbitrary bounded additive noise. That is, each example is of the form  $(x_i, t(x_i) + n_i)$ , where  $t \in \mathcal{H}$  and  $|n_i| < \eta$ . Clearly, any function  $h \in \mathcal{H}$  that is  $\eta$ -close to the training sample will satisfy

$$\Pr(|h - t| \geq 2\eta + \gamma) < \epsilon,$$

provided that  $\mathcal{H}$  approximates from interpolated examples and the training sample is sufficiently large. In addition, if there is an algorithm that can learn in the presence of malicious noise (in this sense), then it can certainly learn in the presence of uniformly distributed random noise (as defined in [3]), which implies  $\text{fat}_{\mathcal{H}}$  is finite ([3, Theorem 3]). That is, a function class  $\mathcal{H}$  is learnable with malicious noise if and only if  $\text{fat}_{\mathcal{H}}$  is finite.

Acknowledgements

This research was supported in part by the Australian Telecommunications and Electronics Research Board. The work of Martin Anthony is supported in part by the European Union through the ‘Neurocolt’ ESPRIT Working Group. The research reported here was conducted while Martin Anthony was visiting the Department of Systems Engineering, Research School of Information Sciences and Engineering, Australian National University.

## References

- [1] Alon, N., Ben-David, S., Cesa-Bianchi, N. and Haussler, D. (1997) Scale-sensitive dimensions, uniform convergence, and learnability. *J. Assoc. Comput. Mach.* **44** 615–631.
- [2] Anthony, M. and Biggs, N. (1992) *Computational Learning Theory: An Introduction*, Cambridge University Press.
- [3] Bartlett, P. L., Long, P. M. and Williamson, R. C. (1994) Fat-shattering and the learnability of real-valued functions. In *Proc. Seventh Annual ACM Conference on Computational Learning Theory*, ACM Press, New York.
- [4] Ben-David, S., Cesa-Bianchi, N., Haussler, D. and Long, P. (1995) Characterizations of learnability for classes of  $\{0, \dots, n\}$ -valued functions. *J. Comput. System Sci.* **50** 74–86. (An earlier version appeared in *Proc. Fifth Annual ACM Workshop on Computational Learning Theory*, ACM Press, New York.)
- [5] Blumer, A., Ehrenfeucht, A., Haussler, D. and Warmuth, M. K. (1989) Learnability and the Vapnik–Chervonenkis dimension. *J. Assoc. Comput. Mach.* **36** 929–965.
- [6] Dudley, R. M., Giné, E. and Zinn, J. (1991) Uniform and universal Glivenko–Cantelli classes. *J. Theoret. Probab.* **4** 485–510.
- [7] Ehrenfeucht, A., Haussler, D., Kearns, M. and Valiant, L. (1989) A general lower bound on the number of examples needed for learning. *Inform. Comput.* **82** 247–261.
- [8] Haussler, D. (1992) Decision theoretic generalizations of the PAC model for neural net and other learning applications. *Inform. Comput.* **100** 78–150.
- [9] Kearns, M. J. and Schapire, R. E. (1994) Efficient distribution-free learning of probabilistic concepts. *J. Comput. System Sci.* **48** 464–497.
- [10] Kolmogorov, A. N. and Tihomirov, V. M. (1961)  $\epsilon$ -entropy and  $\epsilon$ -capacity of sets in functional spaces. *AMS Translations Ser. 2* **17** 277–364.
- [11] Natarajan, B. K. (1993) Occam’s razor for functions. In *Proc. Sixth Annual Workshop on Computational Learning Theory*, ACM Press, New York, pp. 370–376.
- [12] Pollard, D. (1984) *Convergence of Stochastic Processes*, Springer.
- [13] Vapnik, V. N. and Chervonenkis, A. Ya. (1971) On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications* **16** 264–280.