

SOME RESULTS IN STATISTICAL LEARNING THEORY WITH RELEVANCE TO NONLINEAR SYSTEM IDENTIFICATION

Robert C. Williamson

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

Abstract: Statistical Learning Theory comprises a collection of techniques that have been developed in order to theoretically analyse the performance of neural network and other “learning” algorithms. In this paper, a number of recent results in statistical learning theory are summarised in the context of nonlinear system identification. A top-down approach to the problem is taken, leading to the statement of a number of characterisation results. Specific topics covered include empirical risk minimisation, various types of Glivenko-Cantelli classes, scale-sensitive dimensions, sample complexity and the application to dynamic system identification.

Keywords: learning system, learning algorithm, identification, statistical analysis, models

1. INDUCTION OF MODELS

*The players . . . threw these abstract
formulas at one another displaying the
sequences and possibilities of their science.*
— Herman Hesse: *The Glass Bead Game*

It is well to consider the philosophical foundations of anything one does; that is particularly true when one is trying to solve technologically a philosophically impossible problem. System identification is the induction of models from observed events (data). In general it is impossible in the sense that you can never *know* you have identified the truth: “for all is but a woven web of guesses.” You can not even assert with high probability that you are close to the truth Popper [1981, 1980]. The (implicit) philosophical stance of an engineer identifying a dynamical system might be described as a game. The game is this: *If* certain assumptions I am willing to make (with no epistemological justification) are true, *then* on the basis of data available, how reliable a conclusion can I draw?

It is important to realize that in solving a system identification problem, one is not necessarily seeking to discover the “true” model; one is after a model that

attains a certain level of performance. One could say that of the two goals of a scientific theory (predictive performance and explanatory power, the former is valued rather more). The postulation of a “true” model is sometimes useful for theoretical analysis, but as we shall show, an insightful analysis can be performed without the assumption. Furthermore we shall show how one can theoretically determine the number of data points required to attain a pre-specified level of performance. Nevertheless, one is still attempting to reason inductively, and cognisance needs to be taken of the dangers inherent in doing so.

Black-box identification [Sjöberg et al., 1995, Juditsky et al., 1995] is an approach taken when the engineer neither can hypothesise a physical model nor has a specific parametric model imposed upon him (perhaps from extrinsic considerations). However a model structure which “belongs to families that are known to have good flexibility and have been successful in the past” [Sjöberg et al., 1995, p. 1692] may be sought. This is, of course, inductivism by another name, and appeal to such an infinite regress does not really help [see Popper, 1980, page 29].

The extent to which past success will help is indeed uncertain; however all is not lost. One can derive very

¹ Thanks to Erik Weyer for helpful comments on a draft. This work was supported by the Australian Research Council.

useful insight from an analytical consideration of the game mentioned above. One could say that the goal is to find how black is the box, and what determines its blackness? The game might be called “Assumption Engineering.”

One can play the game in many ways. The best games follow an intrinsic logic rather than an extrinsic one. Thus we might care to determine what are the essential limitations of particular approaches to system identification. For example, to what extent is the *parametrisation* chosen for a class of models an essential feature, and to what extent is it a choice of convenience?

In this paper I will review some recent work in the field of statistical learning theory that relates to these questions. Specifically, I will explain what is now known about the characterisation of problem difficulty in a number of formal frameworks (or specific game rules to stretch our metaphor). I will show by example the extent to which the exact rules can significantly affect the conclusions one can infer. I have attempted to follow the *logical* rather than *chronological* structure of the field in order to argue where to go next. I have included technical definitions of the key ideas discussed in order to make the paper reasonably self-contained, but the reader should not be blinded by those technicalities: above anything else, it is the *flavour* of the results and their logical structure that I wish to communicate. There is a longer version of this paper with material on a number of other related topics.

2. REGRESSION AND NONLINEAR BLACK-BOX STRUCTURES

We will adopt the general setting and some of the notation in Sjöberg et al. [1995], Juditsky et al. [1995] to which the reader is referred for a rather fuller introduction. Suppose we have a black-box dynamical system with input $u(t)$ and output $y(t)$. We observe

$$\begin{aligned} u^t &:= [u(1), \dots, u(t)] \\ y^t &:= [y(1), \dots, y(t)] \end{aligned}$$

and would like to infer a relationship between past observations $[u^{t-1}, y^{t-1}]$ and future outputs $y(t)$. A general model for this is

$$y(t) = g(u^{t-1}, y^{t-1}) + v(t).$$

With this setup, the problem becomes: Choose $g \in \mathcal{G}$, a class of candidate models. The quality of a model is often measured by means of a criterion such as

$$\sum_t \|y(t) - g(u^{t-1}, y^{t-1})\|^2 \quad (1)$$

where the sum is over input-output pairs not used in choosing g . Writing $\phi(t) := \phi(u^{t-1}, y^{t-1})$ for the *regression vector* (some *fixed* function of the data), a more precise statement of the problem becomes: Choose $\phi(\cdot)$, choose \mathcal{G} , then pick $g \in \mathcal{G}$ which gives

the “best” model according to a criterion such as (1). However this becomes a general problem of statistical regression. Thus to a large extent, the study of nonlinear black-box models for system identification is the study of nonlinear regression, perhaps with some added complications due to the dependence structure inherent in $\phi(u^{t-1}, y^{t-1})$.

3. A GENERAL FRAMEWORK FOR REGRESSION PROBLEMS

Each discipline which seized upon the Game created its own language of formulas, abbreviations, and possible combinations.

We would like to theoretically understand the intrinsic limitations of the above regression problem. To that end we want to be very careful in proceeding with an analysis that we do not infer fundamental limitations of our analysis *tools*. One recipe to avoid doing this is to aim for *characterisations* of the problem in hand: that is necessary and sufficient conditions to achieve the goals we set ourselves. We can then be reasonably sure that the only place we have ourselves imposed an answer on ourselves is in the *initial* formulation of the problem. This is the approach adopted by Vapnik and we borrow heavily from his books Vapnik [1995, 1982]. A concise presentation of a number of the ideas below can also be found in Vapnik [1993]. We have omitted without further mention when measurability assumptions must be made. (See van der Vaart and Wellner [1996] for a detailed discussion.)

We require a *loss function*. Let \mathcal{X} be the “input” space and \mathcal{Y} the “output” space. Then a *loss function* is a map $L : \mathcal{Y} \times \mathcal{Y}$ which is used to assign a cost of a particular hypothesis $f : \mathcal{X} \rightarrow \mathcal{Y}$ on a given *data point* $(x, y) \in \mathcal{X} \times \mathcal{Y}$ via $\phi((x, y), f) := L(y, f(x))$. We judge the overall performance of f by the expected value of the loss function (or the *risk*) where the expectation is taken with respect to $P_{\mathcal{X}}$, the distribution of the x data points:

$$R(f) := \mathbf{E}[L(Y, f(X))] = \int \phi(z, f) dP(z)$$

where we have written $z = (x, y)$, $\phi(z, f) := L(y, f(x))$ and P is the distribution of (X, Y) on $\mathcal{X} \times \mathcal{Y}$. We can again restate our problem: Given $(x_i, y_i)_{i=1}^n$, $y_i = f^*(x_i) + e_i$, (e_i is a “noise” sequence) and a class of functions \mathcal{F} , determine

$$\hat{f}^* := \arg \min_{f \in \mathcal{F}} R(f).$$

This is impossible, as we are only given a *sample* of n data points and thus do not know P and hence can not compute the risk.

We thus need to adopt an *inductive principle*, which (we hope!) will lead to estimates somehow “close” to \hat{f}^* . Perhaps the most widely used inductive principle is called *Empirical Risk Minimisation* (ERM). There one picks an hypothesis from \mathcal{F} via

$$\hat{f}_n := \arg \min_{f \in \mathcal{F}} R_{\text{emp}}(f)$$

where the *empirical risk* is defined by

$$R_{\text{emp}}(f) := \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i))$$

and is clearly computable from the data provided. The “hope” we mentioned can be expressed formally as requiring that the estimator \hat{f}_n satisfies the two relationships

$$\begin{aligned} R(\hat{f}_n) &\longrightarrow \inf_{f \in \mathcal{F}} R(f) \\ R_{\text{emp}}(\hat{f}_n) &\longrightarrow \inf_{f \in \mathcal{F}} R(f) \end{aligned}$$

in probability as $n \rightarrow \infty$. (This is a minimal requirement; we would also like to know how fast the convergence is — more on that below.)

It turns out that a modification of the above requirements is necessary to rule out certain pathologies. One particular modification is given below, although it is not known if it is the “best” one. The situations ruled out include a coded solution whereby a single observation (data point) suffices to identify exactly a function from within some large class — see Bartlett et al. [1996] for an explicit example. As in Vapnik and Chervonenkis [1991], let $\Lambda(c) := \{f \in \mathcal{F} : R(f) > c\}$. Say that ERM is *strictly consistent* if for all c such that $\Lambda(c) \neq \emptyset$, $\inf\{R_{\text{emp}}(f) : f \in \Lambda(c)\} \rightarrow \inf\{R(f) : f \in \Lambda(c)\}$ in probability as $n \rightarrow \infty$.

Theorem 1. Vapnik and Chervonenkis [1991] Suppose there exist constants A and B such that for all $f \in \mathcal{F}$, $A \leq R(f) \leq B$. Then ERM using \mathcal{F} is strictly consistent if and only if $R_{\text{emp}}(f)$ converges uniformly to $R(f)$ in the sense that for all $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} P\{\sup_{f \in \mathcal{F}} (R(f) - R_{\text{emp}}(f)) > \varepsilon\} = 0. \quad (2)$$

A necessary and sufficient condition for *uniform one-sided convergence* (2) is known, but it is too technical for this short paper. We do however describe the necessary and sufficient conditions for *uniform two-sided convergence*

$$\lim_{n \rightarrow \infty} P\{\sup_{f \in \mathcal{F}} |R(f) - R_{\text{emp}}(f)| > \varepsilon\} = 0 \quad (3)$$

for all $\varepsilon > 0$. In order to state the conditions, we need to introduce a swag of notation.

Let A be a set in a metric space (X, ρ) . A set U is an ε -*covering* of A if for all $a \in A$, there is a $u \in U$ with $\rho(a, u) < \varepsilon$. The ε -*covering number* of A (with respect to ρ) $N_\rho^A(\varepsilon)$ is the number of elements in the smallest ε -cover of A . If f and g are functions defined on X and $x_1, \dots, x_n \in X$, then the ℓ_p^n metric induced by the points is defined by $\rho_p(f, g) := (\frac{1}{n} \sum_{i=1}^n |f(x_i) - g(x_i)|^p)^{1/p}$ (for $p \in [1, \infty)$) and $\max\{|f(x_i) - g(x_i)| : i = 1, \dots, n\}$ for $p = \infty$. If A is a class of functions on X we write $N_p^A(\varepsilon)$ for $N_{\rho_p}^A(\varepsilon)$. Let $\Phi := \{\phi = \phi(z, f) = L(y, f(x)) :$

$z = (x, y), f \in \mathcal{F}\}$ be the *loss-function induced class*. Let $N_p^\Phi(\varepsilon; z_1, \dots, z_n)$ denote the ε -covering number of Φ with respect to z_1, \dots, z_n in the ℓ_p metric. Let $H_p^\Phi(\varepsilon; z_1, \dots, z_n) := \ln N_p^\Phi(\varepsilon; z_1, \dots, z_n)$. The quantity $H_\infty^\Phi(\varepsilon, n) := \mathbf{E}H_\infty^\Phi(\varepsilon, z_1, \dots, z_n)$ is known as the *VC-entropy* of Φ . The expectation is taken with $\underline{z} = (z_1, \dots, z_n)$ drawn according to P^n . Vapnik and Chervonenkis [1981] have shown that if for all $\varepsilon > 0$, $\lim_{n \rightarrow \infty} \frac{1}{n} H_1^\Phi(\varepsilon, n) = 0$ then (3) holds. Furthermore if (3) holds, then for all $\varepsilon > 0$, $\lim_{n \rightarrow \infty} \frac{1}{n} H_\infty^\Phi(\varepsilon, n) = 0$. This is the “best” way of stating the result in the sense that from the definitions, it follows that $H_\infty^\Phi(\varepsilon, n) \geq H_1^\Phi(\varepsilon, n)$. (Use of ℓ_1 covering numbers can lead to better finite sample size bounds.) Traditionally the result has been stated in the more compact form:

Theorem 2. Equation 3 holds if and only if $\forall \varepsilon > 0$

$$\lim_{n \rightarrow \infty} \frac{1}{n} H_\infty^\Phi(\varepsilon, n) = 0. \quad (4)$$

It can equally well be stated in terms of $H_1^\Phi(\varepsilon, n)$; [see Vidyasagar, 1997, page 127].

It would seem however that we have made little progress since condition (4) still depends on the (unknown) probability distribution because of the expectation. Furthermore, even if we did know P , the actual calculation of $H_p^\Phi(\varepsilon, n)$ could be very difficult. Both of these difficulties can be overcome by asking a harder question. If instead of requiring uniform convergence just for a particular P , we can require it to hold uniformly for all P . We now set

$$G_\infty^\Phi(\varepsilon, n) := \sup_{z_1, \dots, z_n} H_\infty^\Phi(\varepsilon; z_1, \dots, z_n)$$

which is (in principle) calculable knowing only \mathcal{F} and L . Clearly, for any distribution P , and any $\varepsilon > 0$, we have $H^\Phi(\varepsilon, n) < G^\Phi(\varepsilon, n)$. In fact ERM is consistent for \mathcal{F} for any probability measure if and only if for all $\varepsilon > 0$, $\lim_{n \rightarrow \infty} \frac{1}{n} G_\infty^\Phi(\varepsilon, n) = 0$ where as above Φ is the loss-function class induced from \mathcal{F} . (The sufficiency in this result follows immediately; necessity follows from the construction of an appropriate P .)

4. GLIVENKO-CANTELLI CLASSES

Whilst in principle $G_\infty^\Phi(\varepsilon, n)$ is now calculable, without some further tricks, it is difficult to do so directly from the definition. In order to make the logical structure of the problem clearer and to link in with some of the literature on the topic, we introduce some further terminology. A key idea is the *dimension* of a class of hypotheses \mathcal{F} . That hypothesis classes with larger “dimensions” require more empirical data to work with is an old idea; for example it was quite explicitly discussed over 60 years ago in [Popper, 1980, Sections 38–39].

Suppose P is a probability distribution on the input space \mathcal{X} . Let $\underline{x}^n := (x_1, \dots, x_n) \in \mathcal{X}^n$. Write

$P(f) := \int f(x)dP(x)$ and $P_n(f) := \frac{1}{n} \sum_{i=1}^n f(x_i)$. A class of real-valued functions \mathcal{F} defined on \mathcal{X} is called a *P-Glivenko-Cantelli* (*P-GC*) class [van der Vaart and Wellner, 1996] if for all $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} P^n \{ \underline{x}^n \in \mathcal{X}^n : \sup_{f \in \mathcal{F}} |P(f) - P_n(f)| > \varepsilon \} = 0.$$

Thus the condition (3) is equivalent to Φ being *P-GC*. We say that \mathcal{F} is an ε -uniform GC class (ε -UGC) if

$$\limsup_{n \rightarrow \infty} \sup_P P^n \{ \underline{x}^n \in \mathcal{X}^n : \sup_{f \in \mathcal{F}} |P(f) - P_n(f)| > \varepsilon \} = 0$$

where the outer supremum is taken over all probability measures P on \mathcal{X} . We say \mathcal{F} is a *uniform GC class* if it is ε -UGC for all $\varepsilon > 0$.

4.1 Fat-Shattering Dimensions

We now introduce two so-called scale-sensitive dimensions which we will use below. Let \mathcal{F} be a set of real valued functions. We say that a set of points X is γ -shattered by \mathcal{F} relative to $r = (r_x)_{x \in X}$ if there are real numbers r_x indexed by $x \in X$ such that for all binary vectors b indexed by X , there is a function $f_b \in \mathcal{F}$ satisfying

$$f_b(x) \begin{cases} \geq r_x + \gamma & \text{if } b_x = 1 \\ \leq r_x - \gamma & \text{otherwise.} \end{cases}$$

The *fat shattering dimension* $\text{Fat}_{\mathcal{F}}$ of the set \mathcal{F} is a function from the positive real numbers to the integers which maps a value γ to the size of the largest γ -shattered set, if this is finite, or infinity otherwise.

If X can be γ -shattered choosing $r_x = r$ the same for all $x \in X$, we say the points X are *level γ -shattered* (at level r) and denote the largest number of such points by the *level fat shattering function* $\text{LFat}_{\mathcal{F}}(\gamma)$. Fat and LFat can not be too different as shown in Alon et al. [1997]: If $\mathcal{F} \subseteq [0, 1]^X$, then for all $\gamma > 0$, $\text{LFat}_{\mathcal{F}}(\gamma) \leq \text{Fat}_{\mathcal{F}}(\gamma) \leq (2 \lceil 1/2\gamma \rceil - 1) \text{Fat}_{\mathcal{F}}(\gamma/2)$.

The fat-shattering dimension for a number of function classes has been computed. An easy example is the set $BV \subseteq \mathbf{R}^{[0,1]}$ of functions of bounded variation for which $\text{Fat}_{BV}(\gamma) = O(1/\gamma)$.

The idea of fat-shattering was introduced into the statistical learning theory community in Kearns and Schapire [1994], although it has been used in approximation theory since the late 1950's where it was apparently first proposed by Kolmogorov [see Tikhomirov, 1960, page 103]. The intuitively appealing thing about $\text{Fat}_{\mathcal{F}}(\gamma)$ is that it measures the complexity of \mathcal{F} at the accuracy scale one is interested in working at. This is in contrast to previous generalisations of the ‘‘Vapnik-Chervonenkis dimension’’ (VC-dim) to treat the problem of learning real-valued functions. For example, the quantity known as the Pollard *pseudo-dimension* (see [Haussler, 1992]) $\text{Pdim}(\mathcal{F})$ can be expressed as $\text{Pdim}(\mathcal{F}) = \lim_{\gamma \rightarrow 0} \text{Fat}_{\mathcal{F}}(\gamma)$. There are many function classes for which $\text{Pdim}(\mathcal{F})$ is infinite or very

large whilst $\text{Fat}_{\mathcal{F}}(\gamma)$ is small for all γ of the appropriate scale (which depends on the accuracy to which one wishes to learn; see the sample complexity results in the next Section).

We say a triple (A, γ, r) , $A \subseteq \mathcal{X}$, $P(A) > 0$, $\gamma > 0$, $r \in \mathbf{R}$, is a *witness of irregularity* if 1) $P|_A$ has no atoms (i.e. $\forall x \in A, P(x) = 0$) and 2) for all $n \geq 1$, $P^n \{ \underline{x}^n \in A^n : \mathcal{F} \text{ } \gamma\text{-shatters } \underline{x}^n \text{ at level } r \} = 1$. The second condition can be stated in words as for all $n \geq 1$, almost all subsets of size n are γ -shattered by \mathcal{F} at level r .

4.2 Characterisations via Fat-Shattering Dimensions

Theorem 3. Talagrand [1987, 1996] \mathcal{F} is *P-GC* if and only if there is no witness of irregularity.

In hindsight, consideration of Talagrand's result along with Vapnik and Chervonenkis' characterisation of uniform GC classes in terms of G^Φ might lead one to conjecture that \mathcal{F} is *uniform-GC* iff $\forall \gamma > 0, \forall r \in \mathbf{R}, \exists n', \forall n > n'$, there is no sequence of n points x_1, \dots, x_n that is γ -shattered at level r . That is indeed the case and although not formally derived from Talagrand's result, was proven in Alon et al. [1997] where it is stated in the following more elegant form:

Theorem 4. The following are equivalent:

- 1) \mathcal{F} is uniform-GC.
- 2) $\text{Fat}_{\mathcal{F}}(\gamma) < \infty$ for all $\gamma > 0$.
- 3) $\text{LFat}_{\mathcal{F}}(\gamma) < \infty$ for all $\gamma > 0$.

One of the attractive things about this result is that it solely depends on ‘‘combinatorial’’ properties of \mathcal{F} . Furthermore, as we shall see in Section 5, these dimensions can be used to make detailed statements about the performance of learning algorithms on a finite number of sample points.

A key technical tool in proving these characterisations is the upper bounding of a ℓ_∞^n covering number in terms of $\text{Fat}_{\mathcal{F}}(\gamma)$. We state the result slightly differently to Alon et al. [1997]; the proof that it is equivalent is given in Lee [1996].

Theorem 5. Suppose \mathcal{F} is a class of $[0, 1]$ -valued functions on \mathcal{X} and $0 < \varepsilon \leq 1$. Then for all $\underline{x}^n \in \mathcal{X}^n$,

$$G_\infty^{\mathcal{F}}(\varepsilon, n) \leq 3 \text{Fat}_{\mathcal{F}}(\varepsilon/4) \ln^2 \frac{16n}{\varepsilon^2}$$

(This result plays a role directly analogous to the Sauer-Shelah lemma for the theory of learning $\{0, 1\}$ -valued functions based on VCdim .)

5. SAMPLE COMPLEXITY

The characterisations of various notions of GCness are but a starting point. They actually lead to insightful

bounds on the sample complexity $m(\varepsilon, \delta)$ of learning problems, which is defined as the number of samples m needed to learn to accuracy ε with probability $1 - \delta$. The precise definitions vary according to the exact setting of the learning problem, and we will only consider the two settings known as regression and agnostic learning. The significance of these results is that they give a *finite* bound on the number of samples need to learn to some pre-specified accuracy. Their practical significance is not so much the specific numbers they indicate, but in the scaling laws (in terms of ε) which are valid for “small” sample sizes.

5.1 Agnostic Learning

Whilst the traditional regression framework certainly provides some insight, it is unsatisfying in one important respect: it requires one to assume that the “true” underlying function $f^*(x) = \mathbf{E}[Y|X = x]$ is in the class \mathcal{F} . This is obviously something that can not be verified. However if we relax the assumption we are faced with another problem: there will now be *two* sources of error. As well as the “statistical” error incurred as before due to only having a finite sample of data, there will also be an “approximation” error between the “truth” and the best model in the class. This decomposition is additive in the case of a square loss function, and we will restrict our attention mainly to that case. It is often explained in terms of bias plus variance. Whilst one can not theoretically quantify the approximation error without (unverifiable) assumptions on the truth, it is still possible to make theoretical claims about the sampling error.

The agnostic learning model used here is based on that in Kearns et al. [1994]. We assume as before that we are presented with a set of examples $\{(x_i, y_i) : i = 1, \dots, n\}$ with $y_i \in \mathcal{Y}$ which is a bounded subset of \mathbf{R} . The examples are drawn independently from some arbitrary distribution P on $\mathcal{X} \times \mathcal{Y}$. A class \mathcal{F} of real-valued functions defined on \mathcal{X} is said to be *agnostically learnable* with *sample complexity* $m(\varepsilon, \delta)$ if for any probability distribution P on $\mathcal{X} \times \mathcal{Y}$, given $0 < \delta \leq 1$ and $\varepsilon > 0$, using $m(\varepsilon, \delta)$ examples one can find a hypothesis $h \in \mathcal{F}$ such that with probability at least $1 - \delta$,

$$R(h) \leq \inf_{f \in \mathcal{F}} R(f) + \varepsilon.$$

Theorem 6. Suppose $\mathcal{F} \subseteq [0, 1]^{\mathcal{X}}$, $\mathcal{Y} \subseteq [0, 1]$ and that for some $r > 0$, $\text{Fat}_{\mathcal{F}}(\gamma) = O(\gamma^{-r})$. Then the sample complexity of agnostically learning with \mathcal{F} satisfies

$$m(\varepsilon, \delta) = O\left(\frac{1}{\varepsilon^2} \left(\text{Fat}_{\mathcal{F}}(\varepsilon) \ln^2 \frac{\text{Fat}_{\mathcal{F}}(\varepsilon)}{\varepsilon} + \ln \frac{1}{\delta}\right)\right).$$

See [Lee, 1996] for a proof.

The key thing to note is that the sample complexity for agnostic learning scales roughly as $1/\varepsilon^2$. It can be

shown that the rate is $1/\varepsilon$ for the traditional regression setting. Note that the constants implicit in the $O(\cdot)$ can be explicitly evaluated.

5.2 Application to System Identification

Results along the lines of the ones presented can be applied to the overall system identification task. Rather than catalog the existing results, let us be satisfied with a sampling of the literature, and some remarks on future directions.

To the best of my knowledge the first application of results of the flavour presented above to problems of System Identification were in Weyer [1992], Weyer et al. [1992, 1993, 1996] where results of Vapnik were used in a setting where the *models* were linear systems, although the true plant was not necessarily. Since then these sorts of results have been applied by various authors (including Meir, Fiechter, Campi and Kumar). One difficulty in applying the above (or related) results directly is the independence assumption required for the above results to hold. This was overcome in Weyer et al. [1996] by making assumptions about the length of the tail of the impulse response of the unknown system (which does not have to be linear). There is now a growing body of work extending the characterisation results (and hence opening the way for sample complexity results) to dependent processes. One example is Peškir and Yukich [1994]. The usual methods of theoretically dealing with dependent processes (by reducing the problem to an approximately independent process, with fewer observations) have been used.

The *point* of applying the general techniques to System Identification problems is severalfold. The key ones are 1) *Finite* sample complexity bounds can be obtained in a probabilistic setting where all of the traditional results are asymptotic; 2) The results tend to show more clearly what it is about a problem that makes it hard; 3) The further development of the general framework, and in particular the possibility of new inductive principles, could provide a range of new algorithms for System Identification.

6. CONCLUSIONS

A number of results from statistical learning theory have been presented. They are of direct relevance to Nonlinear System Identification in so far as the general problem formulation adopted from Sjöberg et al. [1995], Juditsky et al. [1995] is valid, for once that setup is adopted, one is faced with a general regression problem. We have chosen to present the logical structure of the state of knowledge of the field, especially concentrating on the ERM inductive principle. An important conclusion is that the number of parameters in a model class is *not* a characterization

of problem difficulty. This follows from the fact that as we have seen, $\text{Fat}_{\mathcal{F}}(\gamma)$ is such a characterization, and one can readily construct a variety of models with the same numbers of parameters but vastly differing $\text{Fat}_{\mathcal{F}}(\gamma)$. We have omitted mention of some key developments arising in the field within the last couple of years. Two very exciting ones are the rigorous analysis of alternate inductive principles [Shawe-Taylor et al., 1998] and the direct calculation of the relevant covering numbers *without* recourse to combinatorial dimensions [Williamson et al., 1998].

References

- N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler. Scale sensitive dimensions, uniform convergence and learnability. *Journal of the ACM*, 44(4):615–631, 1997.
- P.L. Bartlett, P.M. Long, and R.C. Williamson. Fat-shattering and the learnability of real-valued functions. *Journal of Computer and System Sciences*, 52(3):434–452, 1996.
- D. Haussler. Decision theoretic generalizations of the PAC model for neural net and other learning applications. *Inform. Comput.*, 100(1):78–150, September 1992.
- A. Juditsky, H. Hjalmarsson, A. Benveniste, B. Delyon, L. Ljung, J. Sjöberg, and Q. Zhang. Nonlinear black-box models in system identification: Mathematical foundations. *Automatica*, 31(12):1725–1750, 1995.
- M. J. Kearns, R. E. Schapire, and L. M. Sellie. Toward efficient agnostic learning. *Machine Learning*, 17(2):115–141, 1994.
- M.J. Kearns and R.E. Schapire. Efficient distribution-free learning of probabilistic concepts. *Journal of Computer and System Sciences*, 48(3):464–497, 1994.
- W.S. Lee. *Agnostic Learning and Single Hidden Layer Neural Networks*. PhD thesis, Australian National University, 1996. <http://routh.ee.adfa.oz.au/~weesun/thesis.ps.Z>.
- G. Peškir and J.E. Yukich. Uniform ergodic theorems for dynamical systems under vc entropy conditions. In *Probability in Banach Spaces 9*, pages 105–128. Birkhäuser, Boston, 1994.
- K.R. Popper. *The Logic of Scientific Discovery*. Hutchinson, London, 1980. Originally published in 1934.
- K.R. Popper. *Realism and the Aim of Science*. Rowman and Littlefield, Totowa, New Jersey, 1981.
- J. Shawe-Taylor, P.L. Bartlett, R.C. Williamson, and M. Anthony. Structural risk minimization over data-dependent hierarchies. *IEEE Transactions on Information Theory*, 1998. to appear.
- J. Sjöberg, Q. Zhang, L. Ljung, A. Benveniste, B. Delyon, P.-Y. Glorennec, H. Hjalmarsson, and A. Juditsky. Nonlinear black-box models in system identification: Mathematical foundations. *Automatica*, 31(12):1691–1724, 1995.
- M. Talagrand. The Glivenko-Cantelli problem. *Annals of Probability*, 6:837–870, 1987.
- Michel Talagrand. The Glivenko-Cantelli problem, ten years later. *Journal of Theoretical Probability*, 9(2):371–384, 1996.
- V.M. Tikhomirov. Diameters of sets in function spaces and the theory of best approximations. *Russian Mathematical Surveys*, 15(3):75–111, 1960.
- A.W. van der Vaart and J.A. Wellner. *Weak Convergence and Empirical Processes*. Springer, 1996.
- V.N. Vapnik. *Estimation of Dependences from Empirical Data*. Springer, New York, 1982.
- V.N. Vapnik. Three fundamental concepts of the capacity of learning machines. *Physica A*, 200:838–844, 1993.
- V.N. Vapnik. *The Nature of Statistical Learning Theory*. Springer, New York, 1995.
- V.N. Vapnik and A.Ya. Chervonenkis. Necessary and sufficient conditions for the uniform convergence of means to their expectations. *Theory of Probability and its Applications*, 26(3):532–553, 1981.
- V.N. Vapnik and A.Ya. Chervonenkis. The necessary and sufficient conditions for consistency in the empirical risk method. *Pattern Recognition and Image Analysis*, 1(3):283–305, 1991.
- M. Vidyasagar. *A Theory of Learning and Generalization*. Springer, London, 1997.
- E. Weyer. *System Identification in the Behavioural Framework*. PhD thesis, The Norwegian Institute of Technology, 1992.
- E. Weyer, I.M.Y. Mareels, and R.C. Williamson. Sample complexity of least squares identification of FIR and ARX models. In *Proceedings of the 13th IFAC World Congress*, volume J, pages 239–244, 1996.
- E. Weyer, R.C. Williamson, and I.M.Y. Mareels. An approach to system identification based on risk minimization and behaviours. In *Proceedings of the 31st Conference on Decision and Control*, pages 927–932. IEEE Press, 1992.
- E. Weyer, R.C. Williamson, and I.M.Y. Mareels. A principle for system identification in the behavioural framework. In *Proceedings of the 12th World Congress of the International Federation of Automatic Control*, pages 387–390, 1993.
- R.C. Williamson, A. Smola, and B. Schölkopf. Entropy numbers, operators and support vector kernels. Typescript, 1998.