

Predictive PAC learnability: a paradigm for learning from exchangeable input data

Vladimir Pestov

Department of Mathematics and Statistics

University of Ottawa

Ottawa, Ontario, Canada

vpest283@uottawa.ca

Abstract—Exchangeable random variables form an important and well-studied generalization of i.i.d. variables, however simple examples show that no nontrivial concept or function classes are PAC learnable under general exchangeable data inputs X_1, X_2, \dots . Inspired by the work of Berti and Rigo on a Glivenko–Cantelli theorem for exchangeable inputs, we propose a new paradigm, adequate for learning from exchangeable data: predictive PAC learnability. A learning rule \mathcal{L} for a function class \mathcal{F} is predictive PAC if for every $\varepsilon, \delta > 0$ and each function $f \in \mathcal{F}$, whenever $|\sigma| \geq s(\delta, \varepsilon)$, we have with confidence $1 - \delta$ that the expected difference between $f(X_{n+1})$ and the image of $f|\sigma$ under \mathcal{L} does not exceed ε conditionally on X_1, X_2, \dots, X_n . Thus, instead of learning the function f as such, we are learning to a given accuracy ε the predictive behaviour of f at the future points $X_i(\omega)$, $i > n$ of the sample path. Using de Finetti’s theorem, we show that if a universally separable function class \mathcal{F} is distribution-free PAC learnable under i.i.d. inputs, then it is distribution-free predictive PAC learnable under exchangeable inputs, with a slightly worse sample complexity.

Index Terms—Exchangeable random variables, de Finetti theorem, predictive PAC learnability.

I. INTRODUCTION

In the classical theory of statistical learning as initiated in [15], [4] (see [14] for a historical and philosophical perspective) data inputs are traditionally modelled by a sequence of i.i.d. random variables (X_i) . Generalizing this approach usually involves easing the i.i.d. restriction on the sequence of inputs, all the while trying to obtain the same conclusions as in the classical theory, namely the uniform convergence of empirical means and subsequently the PAC learnability of a concept or a function class under the usual combinatorial restrictions in terms of shattering. For instance, the i.i.d. condition can be relaxed to that of being an ergodic stationary sequence ([12], p. 9), or a β -mixing sequence [16]. As to α -mixing sequences, they are known to result in the same PAC learnable function classes under a single distribution [17], although it is still unknown whether uniform convergence of empirical means takes place [18]. An interesting recent investigation is [11].

However, at some point this approach hits a wall. Among the best studied classes of dependent stationary random variables are exchangeable random variables [6]; [3], p. 473; [9], [10]. A sequence of r.v. (X_i) is *exchangeable*, if for every finite sequence (i_1, i_2, \dots, i_n) of integers the joint distributions of $(X_{i_1}, X_{i_2}, \dots, X_{i_n})$ and of (X_1, X_2, \dots, X_n) are the same.

According to the famous De Finetti theorem [6], [7], a sequence (X_i) is exchangeable if and only if the joint distribution P on Ω^∞ is a mixture of product distributions (that is, (X_i) is a mixture of a family of i.i.d. random sequences).

A nice illustration and the most extreme example of an exchangeable sequence which is not i.i.d is a sequence of identical copies of one and the same random variable, $X_i = X$, $i = 1, 2, \dots$. The joint distribution of this process is a measure supported on the diagonal of the infinite product space Ω^∞ , which is clearly a mixture of infinite powers of all Dirac point masses on Ω .

Now, it is immediately clear that no nontrivial function class \mathcal{F} on a domain Ω will be PAC learnable under such a data input process: almost every sample path \bar{x} will be constant, $\bar{x} = (x, x, x, \dots)$, thus revealing no information about the values of a function $f \in \mathcal{F}$ away from x . Consequently, if we want to be able to learn from exchangeable data inputs, the paradigm of learnability itself has to be re-examined.

A way out was shown by Berti and Rigo in their visionary note [2] where they prove that the classical Glivenko–Cantelli theorem holds for a sequence (X_i) of exchangeable random variables if and only if the sequence is i.i.d. At the same time, they observe that the classical GC theorem is formally equivalent to the statement about the predictive distribution being approximated by the observed frequency:

$$\sup_t |F_n(t, \omega) - P(X_{n+1} \leq t | X_1, \dots, X_n)(\omega)| \rightarrow 0 \text{ a.s.}$$

Here $F_n(t, \omega) = (1/n) \sum_{i=1}^n I_{(-\infty, t]}(X_i)$ is the empirical mean of the indicator function, and $P(\cdot | X_1, \dots, X_n)$ is the conditional probability. As shown in [2], in this form the statement remains valid if the r.v. (X_i) are exchangeable, and the result can be considered as a conditional (or: predictive) version of the classical Glivenko–Cantelli theorem.

Since the uniform Glivenko–Cantelli theorems are at the heart of statistical learning, one would think that the approach of Berti and Rigo should have consequences for learning from exchangeable inputs. We show that this is indeed the case: by replacing PAC learnability with *predictive* PAC learnability, one arrives at a new broad paradigm of learnability suited for learning under exchangeable inputs.

Say that a function class \mathcal{F} is *predictively PAC learnable* under a given class \mathcal{P} of exchangeable random processes (X_n)

if there exists a *predictive PAC learning rule* for \mathcal{F} under \mathcal{P} , that is, a map \mathcal{L} from the sample space \mathcal{S} to a hypothesis class \mathcal{H} such that

$$P\{\sigma: \mathbb{E}(|(\mathcal{L}(f|_\sigma) - f)(X_{n+1})| | X_1, X_2, \dots, X_n) > \varepsilon\} \rightarrow 0$$

uniformly in $f \in \mathcal{F}$ and $(X_i) \in \mathcal{P}$. This is different from PAC learnability in that the expected value of $|\mathcal{L}(f|_\sigma) - f|$ is replaced with the conditional expectation given X_1, X_2, \dots, X_n . If in particular (X_i) are i.i.d., the above definition is a reformulation of PAC learnability under the family of corresponding laws on the domain Ω .

We show that if a function class \mathcal{F} is distribution-free PAC learnable under the usual assumption that the data sample inputs are i.i.d., then \mathcal{F} is predictively PAC learnable under the class of all sequences of exchangeable data inputs. Our results are obtained under the assumption that \mathcal{F} is universally separable.

II. SETTING FOR LEARNABILITY

Here we review the PAC learnability model [1], [4], [13], [16] in order to fix a precise setting. The *domain*, or *instance space*, $\Omega = (\Omega, \mathcal{A})$ is a *measurable space*, that is, a set Ω equipped with a sigma-algebra of subsets \mathcal{A} . We will assume that Ω is a *standard Borel space*, that is, a complete separable metric space equipped with the sigma-algebra of Borel subsets. For instance, without loss in generality one can always assume that $\Omega = \mathbb{R}^k$ is the Euclidean space.

Denote by $\mathcal{B}(\Omega, [0, 1])$ the collection of all Borel measurable functions from Ω to $[0, 1]$. A *function class* \mathcal{F} is a subfamily of $\mathcal{B}(\Omega, [0, 1])$.

The family $P(\Omega)$ of all probability measures on (Ω, \mathcal{A}) is itself a measurable space, whose sigma-algebra is generated by the functions $\nu \mapsto \nu(A)$ from $P(\Omega)$ to \mathbb{R} , as A runs over \mathcal{A} .

In the PAC learning model, a set \mathcal{P} of probability measures on Ω is fixed. Usually either $\mathcal{P} = P(\Omega)$ is the set of all probability measures (*distribution-free learning*), or $\mathcal{P} = \{\mu\}$ is a single measure (*learning under a fixed distribution*).

A *learning sample* is a pair s consisting of a finite subset σ of Ω and of a function on σ . It is convenient to assume that elements $x_1, x_2, \dots, x_n \in \sigma$ are ordered, and thus the set of all samples (σ, τ) with $|\sigma| = n$ can be identified with $(\Omega \times [0, 1])^n$. For $\sigma \in \Omega^n$ and a function $f \in \mathcal{F}$ we will denote $f \upharpoonright \sigma$ the sample obtained by restricting f to σ .

A *learning rule* is a mapping

$$\mathcal{L}: \bigcup_{n=1}^{\infty} \Omega^n \times [0, 1]^n \rightarrow \mathcal{B}(\Omega, [0, 1]),$$

which is measurable with regard to every Borel structure induced on $\mathcal{B}(\Omega, [0, 1])$ by the distances $L^1(\mu)$, $\mu \in \mathcal{P}$.

A learning rule \mathcal{L} is *consistent* if for every $f \in \mathcal{F}$ and each $\sigma \in \Omega^n$ one has

$$\mathcal{L}(f \upharpoonright \sigma) \upharpoonright \sigma = f \upharpoonright \sigma.$$

Consistent learning rules exist for every function class \mathcal{F} under mild measurability restrictions.

A learning rule \mathcal{L} is *probably approximately correct (PAC)* for the function class \mathcal{F} under the class of measures \mathcal{P} if for every $\varepsilon > 0$

$$\sup_{\mu \in \mathcal{P}} \sup_{f \in \mathcal{F}} P\{\sigma \in \Omega^n: \mathbb{E}_\mu |\mathcal{L}(f \upharpoonright \sigma) - f| > \varepsilon\} \rightarrow 0$$

as $n \rightarrow \infty$. Here P stands for $\mu^{\otimes n}$.

Equivalently, there is a function $s(\varepsilon, \delta)$ (*sample complexity* of \mathcal{L}) such that for each $f \in \mathcal{F}$ and every $\mu \in \mathcal{P}$ an i.i.d. sample σ with $\geq s(\varepsilon, \delta)$ points has the property $\mathbb{E}_\mu |f - \mathcal{L}(f \upharpoonright \sigma)| < \varepsilon$ with confidence $\geq 1 - \delta$.

A function class \mathcal{F} is *PAC learnable under \mathcal{P}* , if there exists a PAC learning rule for \mathcal{F} under \mathcal{P} .

If $\mathcal{P} = P(\Omega)$ is the set of all probability measures, then \mathcal{F} is said to be (distribution-free) *PAC learnable*. At the same time, learnability under intermediate families of measures on Ω has received considerable attention, cf. Chapter 7 in [16].

A closely related concept to that of a PAC learnable class is that of a *uniform Glivenko–Cantelli* function class, that is, a function class \mathcal{F} such that for each $\delta, \varepsilon > 0$ one has, whenever $n \geq s(\delta, \varepsilon)$,

$$\sup_{\mu \in P(\Omega)} P \left\{ \sup_{f \in \mathcal{F}} \left| \mathbb{E}_\mu(f) - \frac{1}{n} S_n(f) \right| \geq \varepsilon \right\} < \delta.$$

One also says that \mathcal{C} has the property of *uniform convergence of empirical means (UCEM property)*. Here $s(\delta, \varepsilon)$ is the *sample complexity* of the uniform Glivenko–Cantelli class (which in general has to be distinguished from the sample complexity of a learning rule).

Every uniform Glivenko–Cantelli function class is PAC learnable, for instance, every consistent learning rule for \mathcal{F} is PAC, with the same learning sample complexity. For concept classes, the converse is also true, though not for function classes in general.

A function class \mathcal{F} is *universally separable* [12] if it contains a countable subfamily \mathcal{F}' with the property that every $f \in \mathcal{F}$ is a pointwise limit of a sequence (f_n) of functions from \mathcal{F}' : for each $x \in \Omega$, one has $f_n(x) \rightarrow f(x)$ as $n \rightarrow \infty$.

Notice that in this paper, we only talk of *potential* learnability, adopting a purely information-theoretic viewpoint.

III. EXCHANGEABLE VARIABLES AND DE FINETTI'S THEOREM

De Finetti's theorem, in its classical form ([6], Ch. IV; [7], Th. 7.2) states that a sequence (X_i) of random variables taking values in a standard Borel space Ω is exchangeable if and only if the joint distribution P of the sequence is a mixture of i.i.d. distributions. More precisely, there exists a probability measure η on the Borel space $P(\Omega)$ of probability measures on Ω (the directing measure) so that

$$P = \int_{P(\Omega)} \theta^\infty \eta(d\theta), \quad (1)$$

in the sense that for every measurable function f on Ω^∞ one has

$$\mathbb{E}(f) = \int \mathbb{E}_{\theta^\infty}(f) \eta(d\theta).$$

In this spirit, θ will denote a (random) element of $P(\Omega)$, and “almost all θ ” is to be understood in the sense of directing measure η .

A slightly different viewpoint, adopted in [9], is to fix a *random measure* ν , that is, a measurable mapping from the basic probability space to $P(\Omega)$. Under this approach, de Finetti’s theorem can be put in the following, essentially equivalent, form. Denote by \mathcal{T} the tail sigma-field on Ω^∞ . Then, conditionally on \mathcal{T} , the sequence (X_i) is i.i.d.:

$$P(\omega \in \cdot \mid \mathcal{T}) = \nu^\infty \text{ a.s.}$$

Note that if $\theta \neq \zeta$, then θ^∞ and ζ^∞ are mutually singular. This follows from a remark of Kakutani [8], p. 223: fix f with $\mathbb{E}_\theta(f) \neq \mathbb{E}_\zeta(f)$, then the empirical mean

$$\frac{1}{n} S_n(f) = \frac{1}{n} \sum_{i=1}^n f(X_i)$$

converges at the same time θ^∞ -a.s. to $\mathbb{E}_\theta(f)$ and ζ^∞ -a.s. to $\mathbb{E}_\zeta(f)$. This observation helps to understand the decomposition (1).

The strong law of large numbers for exchangeable variables (cf. e.g. [10], Eq. (2.2) on p. 185, also [9], Proposition 1.4(i)), says that

$$\frac{1}{n} S_n(f) \rightarrow \mathbb{E}(f \mid \mathcal{T}) \quad (2)$$

almost surely. If $P(A) = 1$, then a.s. $\nu(A) = 1$, that is, for almost all θ , one has $\theta(A) = 1$. Thus, the convergence in (2) takes place θ -a.s. for almost all $\theta \in \Theta$. One concludes:

$$\text{For a.e. } \theta, \quad \mathbb{E}(f(X_1) \mid \mathcal{T}) = \mathbb{E}_\theta(f) \text{ } \theta \text{ a.s.} \quad (3)$$

Informally, the conditional expectation $\mathbb{E}(f(X_1) \mid \mathcal{T})$ given the tail sigma-field is viewed by almost every non-random measure θ as a constant function, identically assuming the value $\mathbb{E}_\theta(f)$.

Lemma 3.1: Let X_1, X_2, \dots be a sequence of exchangeable random variables taking values in a standard Borel space Ω . Then for every measurable function f on Ω , for all i and all $j > n$:

$$\mathbb{E}(\mathbb{E}(f(X_i) \mid \mathcal{T}) \mid X_1, \dots, X_n) = \mathbb{E}(f(X_j) \mid X_1, \dots, X_n)$$

a.s., where \mathcal{T} is the tail sigma-field. Consequently, if \mathcal{G} is a countable family of measurable functions, then one has

$$\begin{aligned} \forall f \in \mathcal{G} \quad & \mathbb{E}(\mathbb{E}(f(X_i) \mid \mathcal{T}) \mid X_1, \dots, X_n) \\ & = \mathbb{E}(f(X_j) \mid X_1, \dots, X_n) \end{aligned}$$

almost surely.

Proof: Because of exchangeability, one can assume without loss in generality that $i = 1$ and $j = n + 1$. Now it is enough to establish the result for indicator functions $f = I_A$ of some generating family of Borel subsets $A \subseteq \Omega$, for instance, by identifying Ω with \mathbb{R} and considering the intervals $A = (-\infty, t]$. In this form, the result has been proved in Berti and Rigo [2], where a stronger assertion appears as formula (7) on p. 389. (Their function $F(t, \omega)$ is equal a.s. to $\mathbb{E}(I_{(-\infty, t]}(X_1) \mid \mathcal{T}) = P(X_1 \leq t \mid \mathcal{T})$, which fact follows

from the definition of $F(t, \omega)$ on p. 386, line - 9 as the a.s. limit of $(1/n)S_n(I_{(-\infty, t]})$ and the strong law of large numbers (2)). The second claim is immediate. ■

IV. PREDICTIVE PAC LEARNABILITY

Definition 4.1: Let X_1, X_2, \dots be an exchangeable sequence of random variables with values in a standard Borel space Ω . Denote P the joint distribution on Ω^∞ . We say that a learning rule \mathcal{L} for a function class \mathcal{F} on Ω is *predictively PAC* with sample complexity $s(\delta, \varepsilon)$ (under the sequence (X_i)), if for every $f \in \mathcal{F}$ and each $\varepsilon, \delta > 0$, whenever $n \geq s(\delta, \varepsilon)$, one has

$$P\{\sigma: \mathbb{E}(|(\mathcal{L}(f \upharpoonright \sigma) - f)(X_{n+1})| \mid X_1, X_2, \dots, X_n) > \varepsilon\} < \delta. \quad (4)$$

If \mathcal{P} is a family of sequences of exchangeable random variables, then we say that a function class \mathcal{F} is *predictively PAC learnable under \mathcal{P}* if it admits a learning rule \mathcal{L} that is predictively PAC under every exchangeable sequence $(X_i) \in \mathcal{P}$, with the sample complexity uniformly bounded by some function $s(\delta, \varepsilon)$. Finally, if \mathcal{F} is predictively PAC learnable under the family of all exchangeable sequences (X_i) , we will simply say that \mathcal{F} is *predictively PAC learnable*.

The following theorem is the main result of the article. It allows to deduce predictive PAC learnability from the distribution-free PAC learnability. The proof bypasses a uniform Glivenko–Cantelli theorem for exchangeable variables.

Theorem 4.2: Let \mathcal{F} be a non-trivial universally separable function class on a standard Borel space Ω which is uniform Glivenko–Cantelli (in the classical sense), with the sample complexity $n = s(\delta, \varepsilon)$. Then \mathcal{F} is predictive PAC learnable with the sample complexity $s(\delta\varepsilon, \varepsilon/2)$ under the family of all sequences of Ω -valued exchangeable random variables.

Proof: For every n , let ε_n be the smallest $\varepsilon > 0$ with the property $s(0.5, \varepsilon) \leq n$. Since \mathcal{F} is non-trivial, that is, contains at least two functions, $\varepsilon_n > 0$. Let \mathcal{F}' be a countable dense subfamily of \mathcal{F} such that every $f \in \mathcal{F}$ is a pointwise limit of a sequence of functions from \mathcal{F}' . For every σ , the set of samples $f \upharpoonright \sigma, f \in \mathcal{F}'$ is clearly dense in the set of samples $f \upharpoonright \sigma, f \in \mathcal{F}$. For this reason, using standard selection theorems (e.g. Theorem 5.3.2 in [5]), one can construct a measurable empirical risk minimization learning rule \mathcal{L} on the set of samples

$$S_n(\mathcal{F}) = \{(f \upharpoonright \sigma) : \sigma \in \Omega^n, f \in \mathcal{F}\},$$

taking values in the countable family \mathcal{F}' and such that for every n and each $(\sigma, s) \in S_n(\mathcal{F})$

$$\frac{1}{n} S_n(\mathcal{L}(s) \upharpoonright \sigma - s) < \varepsilon_n.$$

Notice that for every $n \geq s(\delta, \varepsilon)$, whenever $\delta \leq 0.5$, one has $\varepsilon_0 \leq \varepsilon$, and so $\varepsilon + \varepsilon_0 < 2\varepsilon$. For this reason, and taking into account the uniform Glivenko–Cantelli property of \mathcal{F} , for every $\theta \in P(\Omega)$ and each $f \in \mathcal{F}$ one has

$$P\{\mathbb{E}_\theta(\mathcal{L}(f \upharpoonright \sigma) - f) \geq 2\varepsilon\} < \delta. \quad (5)$$

Now let $f \in \mathcal{F}$ and $\varepsilon, \delta > 0$. According to Eq. (3), for a.e. $\theta \in P(\Omega)$ there is a subset $W = W_\theta \subseteq \Omega$ with $\theta(W) = 1$ and such that for every $\omega \in W$ and each $g \in \{f\} \cup \mathcal{F}'$,

$$\mathbb{E}(g|\mathcal{I})(\omega) = \mathbb{E}_\theta(g).$$

Let $\sigma_n(\omega)$ denote, for short, the sequence of values $X_1(\omega), X_2(\omega), \dots, X_n(\omega)$. Define

$$A = \{\omega: \mathbb{E}(|\mathcal{L}(f \upharpoonright \sigma_n(\omega))(X_1) - f(X_1)| | \mathcal{I})(\omega) < 2\varepsilon\}. \quad (6)$$

For a.e. θ , one has, θ -a.s.,

$$A \cap W_\theta = \{\omega: \mathbb{E}_\theta(|\mathcal{L}(f \upharpoonright \sigma_n(\omega)) - f|) < 2\varepsilon\}. \quad (7)$$

According to (5), once $n \geq s(\delta, \varepsilon)$,

$$\theta(A \cap W_\theta) \geq 1 - \delta,$$

and consequently

$$P(A) = \int \theta(A) \eta(d\theta) \geq 1 - \delta.$$

Because of symmetry, we can replace X_1 in the definition (6) of A with X_{n+1} .

Now we are applying Lemma 3.1 to the countable family of functions $\mathcal{G} = \{f\} \cup \{\mathcal{L}(f \upharpoonright \sigma): \sigma \in \Omega^n\}$. Conditioning on X_1, X_2, \dots, X_n amounts to integrating with respect to the conditional distribution $P(d\omega | X_1, X_2, \dots, X_n)$. One must have

$$P\{\omega: P(A^c | X_1, X_2, \dots, X_n)(\omega) \geq 2\varepsilon\} < \delta\varepsilon^{-1}.$$

We conclude:

$$P\{\sigma \in \Omega^n: \mathbb{E}(|\mathcal{L}(\sigma, f | \sigma) - f| | X_1, X_2, \dots, X_n) < 2\varepsilon\} > 1 - \delta\varepsilon^{-1}.$$

Remark 4.3: The proof can be modified so that $\varepsilon/2$ is replaced with $\varepsilon - \gamma_n$ for an arbitrarily sequence $\gamma_n \downarrow 0$. We have only chosen $\varepsilon/2$ for simplicity. On the other hand, the extra factor of ε added to δ does not make much difference, because — unlike the learning precision ε — the confidence parameter δ is well known to be “cheap”.

Corollary 4.4: Let \mathcal{C} be a universally separable concept class on a standard Borel space Ω having finite VC-dimension d . Then \mathcal{C} admits a learning rule which is predictive PAC learnable with regard to any sequence of exchangeable data inputs, with the sample complexity bound

$$s(\delta, \varepsilon) = \max \left\{ \frac{16d}{\varepsilon} \lg \frac{16e}{\varepsilon}, \frac{8}{\varepsilon} \lg \frac{2}{\delta} + \frac{8}{\varepsilon} \lg \frac{1}{\varepsilon} \right\}.$$

The proof follows from Theorem 4.2 and the sample complexity bound for distribution-free PAC learnability ([16], Theorem 7.8),

$$s(\delta, \varepsilon) = \max \left\{ \frac{8d}{\varepsilon} \lg \frac{8e}{\varepsilon}, \frac{4}{\varepsilon} \lg \frac{2}{\delta} \right\}.$$

V. CONCLUSION

Predictable PAC learnability of a function class \mathcal{F} allows to bound, with high confidence, the probability of misclassification of a value of a classifier function $f \in \mathcal{F}$ at any future data sample $X_i(\omega)$, $i \geq n$, given the values of f on a multisample $X_1(\omega), X_2(\omega), \dots, X_n(\omega)$. Under this version of learnability, the function $f \in \mathcal{F}$ cannot be learned in general, it is only its future values that can be predicted with high confidence. For a large number of problems of statistical learning, this is apparently sufficient.

In statistics, exchangeable random variables and de Finetti’s theorem are at the forefront of an ongoing discussion between frequentists and bayesians. (Cf. [3], p. 475.) There is however no need to enter the fray and choose sides, simply because, in Vapnik’s words [13], p. 720,

“Statistical learning theory does not belong to any specific branch of science: It has its own goals, its own paradigm, and its own techniques.

Statisticians (who have their own paradigm) never considered this theory as part of statistics”.

Thus, our new approach can be seen just as an addition to the classical framework of learning theory, possessing its own inner dynamics and putting forward a number of open questions.

Among the most immediate, let us mention the following three, all concerning Theorem 4.2. Can one maintain the initial sample complexity $s(\delta, \varepsilon)$ in the conclusion of the result? Does the theorem hold under less restrictive measurability assumptions on \mathcal{F} than universal separability, for instance, on an assumption that \mathcal{F} is image admissible Souslin ([5], pages 186–187)? Can one conclude that \mathcal{F} is *consistently* predictive PAC learnable, that is, predictive PAC learnable under *every* consistent learning rule \mathcal{L} ?

ACKNOWLEDGMENT

I am indebted to Claus Köstler from whose seminar and conference presentations I have learned about exchangeable random variables and de Finetti’s theorem.

REFERENCES

- [1] Martin Anthony and Peter Bartlett, *Neural network learning: theoretical foundations*. Cambridge University Press, Cambridge, 1999.
- [2] Patrizia Berti and Pietro Rigo, *A Glivenko-Cantelli theorem for exchangeable random variables*, *Statistics & Probability Letters* **32** (1997), 385–391.
- [3] P. Billingsley, *Probability and measure*, 3rd edition. Wiley Series in Probability and Mathematical Statistics. John Wiley & Sons, Inc., New York, 1995.
- [4] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth, *Learnability and the Vapnik-Chervonenkis dimension*, *Journal of the ACM*, **36**(4) (1985), 929–865.
- [5] R.M. Dudley, *Uniform Central Limit Theorems*, *Cambridge Studies in Advanced Mathematics*, **63**, Cambridge University Press, Cambridge, 1999.
- [6] Bruno de Finetti, *La prévision: ses lois logiques, ses sources subjectives*, *Ann. de l’Inst. Henri Poincaré* **7** (1937), 1–68.
- [7] E. Hewitt and L.J. Savage, *Symmetric measures on Cartesian products*, *Trans. Amer. Math. Soc.* **80** (1955), 470–501.
- [8] S. Kakutani, *On equivalence of infinite product measures*, *Ann. of Math.* (2) **49** (1948), 214–224.

- [9] O. Kallenberg, *Probabilistic Symmetries and Invariance Principles*, Probability and its Applications, Springer, New York, 2005.
- [10] J.F.C. Kingman, *Uses of Exchangeability*, Annals of Prob. **6** (1978), 183–197.
- [11] L. Kontorovich, *Measure Concentration of Strongly Mixing Processes with Applications*, PhD thesis, Carnegie-Mellon University, Machine Learning Department, May 2007, 79+vi pp.
- [12] D. Pollard, *Convergence of Stochastic Processes*. Springer-Verlag, New York, 1984.
- [13] V.N. Vapnik, *Statistical Learning Theory*, Wiley, NY, 1998.
- [14] V. Vapnik, *Estimation of Dependences Based on Empirical Data*. Reprint of the 1982 edition. Afterword of 2006: *Empirical inference science*. Information Science and Statistics, Springer, New York, 2006.
- [15] V. N. Vapnik and A. Ya. Chervonenkis, *On the uniform convergence of relative frequencies of events to their probabilities*. Theory Probab. Appl. **16**, issue 2 (1971), 264–280.
- [16] M. Vidyasagar, *Learning and Generalization, with Applications to Neural Networks*, 2nd Ed., Springer-Verlag, 2003.
- [17] M. Vidyasagar, *Convergence of empirical means with alpha-mixing input sequences, and an application to PAC learning*, Proc.44th IEEE Conf. on Decision and Control, and the European Control Conf 2005, pp. 560–565.
- [18] B. Yu, *Rates of convergence of empirical processes for mixing sequences*, Annals of Prob. **22**(1) (1994), 94–116.