

Chapter 3

Transductive PAC-Bayesian learning

3.1. BASIC INEQUALITIES

3.1.1. THE TRANSDUCTIVE SETTING. In this chapter the *observed* sample $(X_i, Y_i)_{i=1}^N$ will be supplemented with a *test* or *shadow* sample $(X_i, Y_i)_{i=N+1}^{(k+1)N}$. This point of view, called *transductive classification*, has been introduced by V. Vapnik. It may be justified in different ways.

On the practical side, one interest of the transductive setting is that it is often a lot easier to collect examples than it is to label them, so that it is not unrealistic to assume that we indeed have two training samples, one labelled and one unlabelled. It also covers the case when a batch of patterns is to be classified and we are allowed to observe the whole batch before issuing the classification.

On the mathematical side, considering a shadow sample proves technically fruitful. Indeed, when introducing the Vapnik–Cervonenkis entropy and Vapnik–Cervonenkis dimension concepts, as well as when dealing with compression schemes, albeit the *inductive* setting is our final concern, the transductive setting is a useful detour. In this second scenario, intermediate technical results involving the shadow sample are integrated with respect to unobserved random variables in a second stage of the proofs.

Let us describe now the changes to be made to previous notation to adapt them to the transductive setting. The distribution \mathbb{P} will be a probability measure on the canonical space $\Omega = (\mathcal{X} \times \mathcal{Y})^{(k+1)N}$, and $(X_i, Y_i)_{i=1}^{(k+1)N}$ will be the canonical process on this space (that is the coordinate process). Unless explicitly mentioned, the parameter k indicating the size of the shadow sample will remain fixed. Assuming the shadow sample size is a multiple of the training sample size is convenient without significantly restricting generality. For a while, we will use a weaker assumption than independence, assuming that \mathbb{P} is *partially exchangeable*, since this is all we need in the proofs.

DEFINITION 3.1.1. For $i = 1, \dots, N$, let $\tau_i : \Omega \rightarrow \Omega$ be defined for any

$\omega = (\omega_j)_{j=1}^{(k+1)N} \in \Omega$ by

$$\begin{cases} \tau_i(\omega)_{i+jN} = \omega_{i+(j-1)N}, & j = 1, \dots, k, \\ \tau_i(\omega)_i = \omega_{i+kN}, \\ \text{and } \tau_i(\omega)_{m+jN} = \omega_{m+jN}, & m \neq i, m = 1, \dots, N, j = 0, \dots, k. \end{cases}$$

Clearly, if we arrange the $(k+1)N$ samples in a $N \times (k+1)$ array, τ_i performs a circular permutation of $k+1$ entries on the i th row, leaving the other rows unchanged. Moreover, all the circular permutations of the i th row have the form τ_i^j , j ranging from 0 to k .

The probability distribution \mathbb{P} is said to be partially exchangeable if for any $i = 1, \dots, N$, $\mathbb{P} \circ \tau_i^{-1} = \mathbb{P}$.

This means equivalently that for any bounded measurable function $h : \Omega \rightarrow \mathbb{R}$, $\mathbb{P}(h \circ \tau_i) = \mathbb{P}(h)$.

In the same way a function h defined on Ω will be said to be partially exchangeable if $h \circ \tau_i = h$ for any $i = 1, \dots, N$. Accordingly a posterior distribution $\rho : \Omega \rightarrow \mathcal{M}_+^1(\Theta, \mathcal{T})$ will be said to be partially exchangeable when $\rho(\omega, A) = \rho[\tau_i(\omega), A]$, for any $\omega \in \Omega$, any $i = 1, \dots, N$ and any $A \in \mathcal{T}$.

For any bounded measurable function h , let us define $T_i(h) = \frac{1}{k+1} \sum_{j=0}^k h \circ \tau_i^j$. Let $T(h) = T_N \circ \dots \circ T_1(h)$. For any partially exchangeable probability distribution \mathbb{P} , and for any bounded measurable function h , $\mathbb{P}[T(h)] = \mathbb{P}(h)$. Let us put

$$\sigma_i(\theta) = \mathbb{1}[f_\theta(X_i) \neq Y_i], \quad \begin{array}{l} \text{indicating the success or failure of } f_\theta \\ \text{to predict } Y_i \text{ from } X_i, \end{array}$$

$$r_1(\theta) = \frac{1}{N} \sum_{i=1}^N \sigma_i(\theta), \quad \begin{array}{l} \text{the empirical error rate of } f_\theta \\ \text{on the observed sample,} \end{array}$$

$$r_2(\theta) = \frac{1}{kN} \sum_{i=N+1}^{(k+1)N} \sigma_i(\theta), \quad \text{the error rate of } f_\theta \text{ on the shadow sample,}$$

$$\bar{r}(\theta) = \frac{r_1(\theta) + kr_2(\theta)}{k+1} = \frac{1}{(k+1)N} \sum_{i=1}^{(k+1)N} \sigma_i(\theta), \quad \begin{array}{l} \text{the global error} \\ \text{rate of } f_\theta, \end{array}$$

$$R_i(\theta) = \mathbb{P}[f_\theta(X_i) \neq Y_i], \quad \begin{array}{l} \text{the expected error} \\ \text{rate of } f_\theta \text{ on the } i\text{th input,} \end{array}$$

$$R(\theta) = \frac{1}{N} \sum_{i=1}^N R_i(\theta) = \mathbb{P}[r_1(\theta)] = \mathbb{P}[r_2(\theta)], \quad \text{the average expected}$$

error rate of f_θ on all inputs.

We will allow for posterior distributions $\rho : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$ depending on the shadow sample. The most interesting ones will anyhow be independent of the shadow labels $Y_{N+1}, \dots, Y_{(k+1)N}$. We will be interested in the conditional expected error rate of the randomized classification rule described by ρ on the shadow sample, given the observed sample, that is, $\mathbb{P}[\rho(r_2)|(X_i, Y_i)_{i=1}^N]$. This is a natural extension of the notion of *generalization error rate*: this is indeed the error rate to be expected when the randomized classification rule described by the posterior distribution ρ is applied to the shadow sample (which should in this case more purposefully be called the test sample).

To see the connection with the previously defined generalization error rate, let us comment on the case when \mathbb{P} is invariant by any permutation of any row, meaning that

$$\mathbb{P}[h(\omega \circ s)] = \mathbb{P}[h(\omega)] \text{ for all } s \in \mathfrak{S}(\{i + jN; j = 0, \dots, k\})$$

and all $i = 1, \dots, N$, where $\mathfrak{S}(A)$ is the set of permutations of A , extended to $\{1, \dots, (k+1)N\}$ so as to be the identity outside of A . In other words, \mathbb{P} is assumed to be invariant under any permutation which keeps the rows unchanged. In this case, if ρ is invariant by any permutation of any row of the shadow sample, meaning that $\rho(\omega \circ s) = \rho(\omega) \in \mathcal{M}_+^1(\Theta)$, $s \in \mathfrak{S}(\{i + jN; j = 1, \dots, k\})$, $i = 1, \dots, N$, then $\mathbb{P}[\rho(r_2)|(X_i, Y_i)_{i=1}^N] = \frac{1}{N} \sum_{i=1}^N \mathbb{P}[\rho(\sigma_{i+N})|(X_i, Y_i)_{i=1}^N]$, meaning that the expectation can be taken on a restricted shadow sample of the same size as the observed sample. If moreover the rows are equidistributed, meaning that their marginal distributions are equal, then

$$\mathbb{P}[\rho(r_2)|(X_i, Y_i)_{i=1}^N] = \mathbb{P}[\rho(\sigma_{N+1})|(X_i, Y_i)_{i=1}^N].$$

This means that under these quite commonly fulfilled assumptions, the expectation can be taken on a single new object to be classified, our study thus covers the case when only one of the patterns from the shadow sample is to be labelled and one is interested in the expected error rate of this single labelling. Of course, in the case when \mathbb{P} is i.i.d. and ρ depends only on the training sample $(X_i, Y_i)_{i=1}^N$, we fall back on the usual criterion of performance $\mathbb{P}[\rho(r_2)|(Z_i)_{i=1}^N] = \rho(R) = \rho(R_1)$.

3.1.2. ABSOLUTE BOUND. Using an obvious factorization, and considering for the moment a fixed value of θ and any partially exchangeable positive real measurable function $\lambda : \Omega \rightarrow \mathbb{R}_+$, we can compute the log-Laplace transform of r_1 under T , which acts like a conditional probability distribution:

$$\begin{aligned} \log \left\{ T[\exp(-\lambda r_1)] \right\} &= \sum_{i=1}^N \log \left\{ T_i \left[\exp\left(-\frac{\lambda}{N} \sigma_i\right) \right] \right\} \\ &\leq N \log \left\{ \frac{1}{N} \sum_{i=1}^N T_i \left[\exp\left(-\frac{\lambda}{N} \sigma_i\right) \right] \right\} = -\lambda \Phi_{\frac{\lambda}{N}}(\bar{r}), \end{aligned}$$

where the function $\Phi_{\frac{\lambda}{N}}$ was defined by equation (1.1, page 2). Remarking that $T \left\{ \exp \left[\lambda \left[\Phi_{\frac{\lambda}{N}}(\bar{r}) - r_1 \right] \right] \right\} = \exp \left[\lambda \Phi_{\frac{\lambda}{N}}(\bar{r}) \right] T[\exp(-\lambda r_1)]$ we obtain

LEMMA 3.1.1. *For any $\theta \in \Theta$ and any partially exchangeable positive real measurable function $\lambda : \Omega \rightarrow \mathbb{R}_+$,*

$$T \left\{ \exp \left[\lambda \left\{ \Phi_{\frac{\lambda}{N}}[\bar{r}(\theta)] - r_1(\theta) \right\} \right] \right\} \leq 1.$$

We deduce from this lemma a result analogous to the inductive case:

THEOREM 3.1.2. *For any partially exchangeable positive real measurable function $\lambda : \Omega \times \Theta \rightarrow \mathbb{R}_+$, for any partially exchangeable posterior distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$,*

$$\mathbb{P} \left\{ \exp \left[\sup_{\rho \in \mathcal{M}_+^1(\Theta)} \rho \left[\lambda \left[\Phi_{\frac{\lambda}{N}}(\bar{r}) - r_1 \right] \right] - \mathcal{K}(\rho, \pi) \right] \right\} \leq 1.$$

The proof is deduced from the previous lemma, using the fact that π is partially exchangeable:

$$\begin{aligned} & \mathbb{P} \left\{ \exp \left[\sup_{\rho \in \mathcal{M}_+^1(\Theta)} \rho \left[\lambda \left[\Phi_{\frac{\lambda}{N}}(\bar{r}) - r_1 \right] \right] - \mathcal{K}(\rho, \pi) \right] \right\} \\ &= \mathbb{P} \left\{ \pi \left\{ \exp \left[\lambda \left[\Phi_{\frac{\lambda}{N}}(\bar{r}) - r_1 \right] \right] \right\} \right\} = \mathbb{P} \left\{ T\pi \left\{ \exp \left[\lambda \left[\Phi_{\frac{\lambda}{N}}(\bar{r}) - r_1 \right] \right] \right\} \right\} \\ &= \mathbb{P} \left\{ \pi \left\{ T \exp \left[\lambda \left[\Phi_{\frac{\lambda}{N}}(\bar{r}) - r_1 \right] \right] \right\} \right\} \leq 1. \end{aligned}$$

3.1.3. RELATIVE BOUNDS. Introducing in the same way

$$\begin{aligned} m'(\theta, \theta') &= \frac{1}{N} \sum_{i=1}^N \left| \mathbb{1}[f_\theta(X_i) \neq Y_i] - \mathbb{1}[f_{\theta'}(X_i) \neq Y_i] \right| \\ \text{and } \bar{m}(\theta, \theta') &= \frac{1}{(k+1)N} \sum_{i=1}^{(k+1)N} \left| \mathbb{1}[f_\theta(X_i) \neq Y_i] - \mathbb{1}[f_{\theta'}(X_i) \neq Y_i] \right|, \end{aligned}$$

we could prove along the same line of reasoning

THEOREM 3.1.3. *For any real parameter λ , any $\tilde{\theta} \in \Theta$, any partially exchangeable posterior distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$,*

$$\begin{aligned} & \mathbb{P} \left\{ \exp \left[\sup_{\rho \in \mathcal{M}_+^1(\Theta)} \lambda \left[\rho \left\{ \Psi_{\frac{\lambda}{N}}[\bar{r}(\cdot) - \bar{r}(\tilde{\theta})], \bar{m}(\cdot, \tilde{\theta}) \right\} \right. \right. \right. \\ & \quad \left. \left. \left. - [\rho(r_1) - r_1(\tilde{\theta})] \right] - \mathcal{K}(\rho, \pi) \right] \right\} \leq 1, \end{aligned}$$

where the function $\Psi_{\frac{\lambda}{N}}$ was defined by equation (1.21, page 35).

THEOREM 3.1.4. *For any real constant γ , for any $\tilde{\theta} \in \Theta$, for any partially exchangeable posterior distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$,*

$$\begin{aligned} & \mathbb{P} \left\{ \exp \left[\sup_{\rho \in \mathcal{M}_+^1(\Theta)} \left\{ -N\rho \left\{ \log \left[1 - \tanh\left(\frac{\gamma}{N}\right) [\bar{r}(\cdot) - \bar{r}(\tilde{\theta})] \right] \right\} \right. \right. \right. \\ & \quad \left. \left. \left. - \gamma [\rho(r_1) - r_1(\tilde{\theta})] - N \log \left[\cosh\left(\frac{\gamma}{N}\right) \right] \rho [m'(\cdot, \tilde{\theta})] - \mathcal{K}(\rho, \pi) \right\} \right] \right\} \leq 1. \end{aligned}$$

This last theorem can be generalized to give

THEOREM 3.1.5. *For any real constant γ , for any partially exchangeable posterior distributions $\pi^1, \pi^2 : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$,*

$$\begin{aligned} & \mathbb{P} \left\{ \exp \left[\sup_{\rho_1, \rho_2 \in \mathcal{M}_+^1(\Theta)} \left\{ -N \log \left\{ 1 - \tanh\left(\frac{\gamma}{N}\right) [\rho_1(\bar{r}) - \rho_2(\bar{r})] \right\} \right. \right. \right. \\ & \quad \left. \left. \left. - \gamma [\rho_1(r_1) - \rho_2(r_1)] - N \log \left[\cosh\left(\frac{\gamma}{N}\right) \right] \rho_1 \otimes \rho_2(m') \right. \right. \right. \\ & \quad \left. \left. \left. - \mathcal{K}(\rho_1, \pi^1) - \mathcal{K}(\rho_2, \pi^2) \right\} \right] \right\} \leq 1. \end{aligned}$$

To conclude this section, we see that the basic theorems of transductive PAC-Bayesian classification have exactly the same form as the basic inequalities of inductive classification, Theorems 1.1.4 (page 4), 1.4.2 (page 35) and 1.4.3 (page 37) with $R(\theta)$ replaced with $\bar{r}(\theta)$, $r(\theta)$ replaced with $r_1(\theta)$ and $M'(\theta, \tilde{\theta})$ replaced with $\bar{m}(\theta, \tilde{\theta})$.

Thus all the results of the first two chapters remain true under the hypotheses of transductive classification, with $R(\theta)$ replaced with $\bar{r}(\theta)$, $r(\theta)$ replaced with $r_1(\theta)$ and $M'(\theta, \tilde{\theta})$ replaced with $\bar{m}(\theta, \tilde{\theta})$.

Consequently, in the case when the unlabelled shadow sample is observed, it is possible to improve on the Vapnik bounds to be discussed hereafter by using an explicit partially exchangeable posterior distribution π and resorting to localized or to relative bounds (in the case at least of unlimited computing resources, which of course may still be unrealistic in many real world situations, and with the caveat, to be recalled in the conclusion of this study, that for small sample sizes and comparatively complex classification models, the improvement may not be so decisive).

Let us notice also that the transductive setting when experimentally available, has the advantage that

$$\begin{aligned} \bar{d}(\theta, \theta') &= \frac{1}{(k+1)N} \sum_{i=1}^{(k+1)N} \mathbb{1}[f_{\theta'}(X_i) \neq f_{\theta}(X_i)] \\ &\geq \bar{m}(\theta, \theta') \geq \bar{r}(\theta) - \bar{r}(\theta'), \quad \theta, \theta' \in \Theta, \end{aligned}$$

is observable in this context, providing an empirical upper bound for the difference $\bar{r}(\hat{\theta}) - \rho(\bar{r})$ for any non-randomized estimator $\hat{\theta}$ and any posterior distribution ρ , namely

$$\bar{r}(\hat{\theta}) \leq \rho(\bar{r}) + \rho[\bar{d}(\cdot, \hat{\theta})].$$

Thus in the setting of transductive statistical experiments, the PAC-Bayesian framework provides fully empirical bounds for the error rate of non-randomized estimators $\hat{\theta} : \Omega \rightarrow \Theta$, even when using a non-atomic prior π (or more generally a non-atomic partially exchangeable posterior distribution π), even when Θ is not a vector space and even when $\theta \mapsto R(\theta)$ cannot be proved to be convex on the support of some useful posterior distribution ρ .

3.2. VAPNIK BOUNDS FOR TRANSDUCTIVE CLASSIFICATION

In this section, we will stick to plain unlocalized non-relative bounds. As we have already mentioned, (and as it was put forward by Vapnik himself in his seminal works), these bounds are not always superseded by the asymptotically better ones when the sample is of small size: they deserve all our attention for this reason. We will start with the general case of a shadow sample of arbitrary size. We will then discuss the case of a shadow sample of equal size to the training set and the case of a fully exchangeable sample distribution, showing how they can be taken advantage of to sharpen inequalities.

3.2.1. WITH A SHADOW SAMPLE OF ARBITRARY SIZE. The great thing with the transductive setting is that we are manipulating only r_1 and \bar{r} which can take only a finite number of values and therefore are piecewise constant on Θ . This makes it possible to derive inequalities that will hold uniformly for any value of the parameter

$\theta \in \Theta$. To this purpose, let us consider for any value $\theta \in \Theta$ of the parameter the subset $\Delta(\theta) \subset \Theta$ of parameters θ' such that the classification rule $f_{\theta'}$ answers the same on the extended sample $(X_i)_{i=1}^{(k+1)N}$ as f_{θ} . Namely, let us put for any $\theta \in \Theta$

$$\Delta(\theta) = \{\theta' \in \Theta; f_{\theta'}(X_i) = f_{\theta}(X_i), i = 1, \dots, (k+1)N\}.$$

We see immediately that $\Delta(\theta)$ is an exchangeable parameter subset on which r_1 and r_2 and therefore also \bar{r} take constant values. Thus for any $\theta \in \Theta$ we may consider the posterior ρ_{θ} defined by

$$\frac{d\rho_{\theta}}{d\pi}(\theta') = \mathbb{1}[\theta' \in \Delta(\theta)]\pi[\Delta(\theta)]^{-1},$$

and use the fact that $\rho_{\theta}(r_1) = r_1(\theta)$ and $\rho_{\theta}(\bar{r}) = \bar{r}(\theta)$, to prove that

LEMMA 3.2.1. *For any partially exchangeable positive real measurable function $\lambda : \Omega \times \Theta \rightarrow \mathbb{R}$ such that*

$$(3.1) \quad \lambda(\omega, \theta') = \lambda(\omega, \theta), \quad \theta \in \Theta, \theta' \in \Delta(\theta), \omega \in \Omega,$$

and any partially exchangeable posterior distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,

$$\Phi_{\frac{\lambda}{N}}[\bar{r}(\theta)] + \frac{\log\{\epsilon\pi[\Delta(\theta)]\}}{\lambda(\theta)} \leq r_1(\theta).$$

We can then remark that for any value of λ independent of ω , the left-hand side of the previous inequality is a partially exchangeable function of $\omega \in \Omega$. Thus this left-hand side is maximized by some partially exchangeable function λ , namely

$$\arg \max_{\lambda} \left\{ \Phi_{\frac{\lambda}{N}}[\bar{r}(\theta)] + \frac{\log\{\epsilon\pi[\Delta(\theta)]\}}{\lambda} \right\}$$

is partially exchangeable as depending only on partially exchangeable quantities. Moreover this choice of $\lambda(\omega, \theta)$ satisfies also condition (3.1) stated in the previous lemma of being constant on $\Delta(\theta)$, proving

LEMMA 3.2.2. *For any partially exchangeable posterior distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$ and any $\lambda \in \mathbb{R}_+$,*

$$\Phi_{\frac{\lambda}{N}}[\bar{r}(\theta)] + \frac{\log\{\epsilon\pi[\Delta(\theta)]\}}{\lambda} \leq r_1(\theta).$$

Writing $\bar{r} = \frac{r_1 + kr_2}{k+1}$ and rearranging terms we obtain

THEOREM 3.2.3. *For any partially exchangeable posterior distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,*

$$r_2(\theta) \leq \frac{k+1}{k} \inf_{\lambda \in \mathbb{R}_+} \frac{1 - \exp\left(-\frac{\lambda}{N}r_1(\theta) + \frac{\log\{\epsilon\pi[\Delta(\theta)]\}}{N}\right)}{1 - \exp\left(-\frac{\lambda}{N}\right)} - \frac{r_1(\theta)}{k}.$$

If we have a set of binary classification rules $\{f_\theta; \theta \in \Theta\}$ whose Vapnik–Cervonenkis dimension is not greater than h , we can choose π such that $\pi[\Delta(\theta)]$ is independent of θ and not less than $\left(\frac{h}{e(k+1)N}\right)^h$, as will be proved further on in Theorem 4.2.2 (page 144).

Another important setting where the complexity term $-\log\{\pi[\Delta(\theta)]\}$ can easily be controlled is the case of *compression schemes*, introduced by Little et al. (1986). It goes as follows: we are given for each labelled sub-sample $(X_i, Y_i)_{i \in J}$, $J \subset \{1, \dots, N\}$, an estimator of the parameter

$$\hat{\theta}[(X_i, Y_i)_{i \in J}] = \hat{\theta}_J, \quad J \subset \{1, \dots, N\}, |J| \leq h,$$

where

$$\hat{\theta} : \bigsqcup_{k=1}^N (\mathcal{X} \times \mathcal{Y})^k \rightarrow \Theta$$

is an exchangeable function providing estimators for sub-samples of arbitrary size. Let us assume that $\hat{\theta}$ is exchangeable, meaning that for any $k = 1, \dots, N$ and any permutation σ of $\{1, \dots, k\}$

$$\hat{\theta}[(x_i, y_i)_{i=1}^k] = \hat{\theta}[(x_{\sigma(i)}, y_{\sigma(i)})_{i=1}^k], \quad (x_i, y_i)_{i=1}^k \in (\mathcal{X} \times \mathcal{Y})^k.$$

In this situation, we can introduce the exchangeable subset

$$\{\hat{\theta}_J; J \subset \{1, \dots, (k+1)N\}, |J| \leq h\} \subset \Theta,$$

which is seen to contain at most

$$\sum_{j=0}^h \binom{(k+1)N}{j} \leq \left(\frac{e(k+1)N}{h}\right)^h$$

classification rules — as will be proved later on in Theorem 4.2.3 (page 144). Note that we had to extend the range of J to all the subsets of the extended sample, although we will use for estimation only those of the training sample, on which the labels are observed. Thus in this case also we can find a partially exchangeable posterior distribution π such that

$$\pi[\Delta(\hat{\theta}_J)] \geq \left(\frac{h}{e(k+1)N}\right)^h.$$

We see that the size of the compression scheme plays the same role in this complexity bound as the Vapnik–Cervonenkis dimension for Vapnik–Cervonenkis classes.

In these two cases of binary classification with Vapnik–Cervonenkis dimension not greater than h and compression schemes depending on a compression set with at most h points, we get a bound of

$$r_2(\theta) \leq \frac{k+1}{k} \inf_{\lambda \in \mathbb{R}_+} \frac{1 - \exp\left(-\frac{\lambda}{N} r_1(\theta) - \frac{h \log\left(\frac{e(k+1)N}{h}\right) - \log(\epsilon)}{N}\right)}{1 - \exp\left(-\frac{\lambda}{N}\right)} - \frac{r_1(\theta)}{k}.$$

Let us make some numerical application: when $N = 1000, h = 10, \epsilon = 0.01$, and $\inf_{\Theta} r_1 = r_1(\hat{\theta}) = 0.2$, we find that $r_2(\hat{\theta}) \leq 0.4093$, for k between 15 and 17, and values of λ equal respectively to 965, 968 and 971. For $k = 1$, we find only $r_2(\hat{\theta}) \leq 0.539$, showing the interest of allowing k to be larger than 1.

3.2.2. WHEN THE SHADOW SAMPLE HAS THE SAME SIZE AS THE TRAINING SAMPLE. In the case when $k = 1$, we can improve Theorem 3.1.2 by taking advantage of the fact that $T_i(\sigma_i)$ can take only 3 values, namely 0, 0.5 and 1. We see thus that $T_i(\sigma_i) - \Phi_{\frac{\lambda}{N}}[T_i(\sigma_i)]$ can take only two values, 0 and $\frac{1}{2} - \Phi_{\frac{\lambda}{N}}(\frac{1}{2})$, because $\Phi_{\frac{\lambda}{N}}(0) = 0$ and $\Phi_{\frac{\lambda}{N}}(1) = 1$. Thus

$$T_i(\sigma_i) - \Phi_{\frac{\lambda}{N}}[T_i(\sigma_i)] = [1 - |1 - 2T_i(\sigma_i)|] \left[\frac{1}{2} - \Phi_{\frac{\lambda}{N}}\left(\frac{1}{2}\right) \right].$$

This shows that in the case when $k = 1$,

$$\begin{aligned} \log \left\{ T[\exp(-\lambda r_1)] \right\} &= -\lambda \bar{r} + \frac{\lambda}{N} \sum_{i=1}^N T_i(\sigma_i) - \Phi_{\frac{\lambda}{N}}[T_i(\sigma_i)] \\ &= -\lambda \bar{r} + \frac{\lambda}{N} \sum_{i=1}^N [1 - |1 - 2T_i(\sigma_i)|] \left[\frac{1}{2} - \Phi_{\frac{\lambda}{N}}\left(\frac{1}{2}\right) \right] \\ &\leq -\lambda \bar{r} + \lambda \left[\frac{1}{2} - \Phi_{\frac{\lambda}{N}}\left(\frac{1}{2}\right) \right] [1 - |1 - 2\bar{r}|]. \end{aligned}$$

Noticing that $\frac{1}{2} - \Phi_{\frac{\lambda}{N}}(\frac{1}{2}) = \frac{N}{\lambda} \log[\cosh(\frac{\lambda}{2N})]$, we obtain

THEOREM 3.2.4. *For any partially exchangeable function $\lambda : \Omega \times \Theta \rightarrow \mathbb{R}_+$, for any partially exchangeable posterior distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$,*

$$\mathbb{P} \left\{ \exp \left[\sup_{\rho \in \mathcal{M}_+^1(\Theta)} \rho \left[\lambda(\bar{r} - r_1) - N \log \left[\cosh\left(\frac{\lambda}{2N}\right) \right] (1 - |1 - 2\bar{r}|) \right] - \mathcal{K}(\rho, \pi) \right] \right\} \leq 1.$$

As a consequence, reasoning as previously, we deduce

THEOREM 3.2.5. *In the case when $k = 1$, for any partially exchangeable posterior distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$ and any $\lambda \in \mathbb{R}_+$,*

$$\bar{r}(\theta) - \frac{N}{\lambda} \log \left[\cosh\left(\frac{\lambda}{2N}\right) \right] (1 - |1 - 2\bar{r}(\theta)|) + \frac{\log \{ \epsilon \pi[\Delta(\theta)] \}}{\lambda} \leq r_1(\theta);$$

and consequently for any $\theta \in \Theta$,

$$r_2(\theta) \leq 2 \inf_{\lambda \in \mathbb{R}_+} \frac{r_1(\theta) - \frac{\log \{ \epsilon \pi[\Delta(\theta)] \}}{\lambda}}{1 - \frac{2N}{\lambda} \log \left[\cosh\left(\frac{\lambda}{2N}\right) \right]} - r_1(\theta).$$

In the case of binary classification using a Vapnik–Cervonenkis class of Vapnik–Cervonenkis dimension not greater than h , we can choose π such that $-\log \{ \pi[\Delta(\theta)] \} \leq h \log(\frac{2\epsilon N}{h})$ and obtain the following numerical illustration of

this theorem: for $N = 1000$, $h = 10$, $\epsilon = 0.01$ and $\inf_{\Theta} r_1 = r_1(\hat{\theta}) = 0.2$, we find an upper bound $r_2(\hat{\theta}) \leq 0.5033$, which improves on Theorem 3.2.3 but still is not under the significance level $\frac{1}{2}$ (achieved by blind random classification). This indicates that considering shadow samples of arbitrary sizes some noisy situations yields a significant improvement on bounds obtained with a shadow sample of the same size as the training sample.

3.2.3. WHEN MOREOVER THE DISTRIBUTION OF THE AUGMENTED SAMPLE IS EXCHANGEABLE. When $k = 1$ and \mathbb{P} is exchangeable meaning that for any bounded measurable function $h : \Omega \rightarrow \mathbb{R}$ and any permutation $s \in \mathfrak{S}(\{1, \dots, 2N\})$ $\mathbb{P}[h(\omega \circ s)] = \mathbb{P}[h(\omega)]$, then we can still improve the bound as follows. Let

$$T'(h) = \frac{1}{N!} \sum_{s \in \mathfrak{S}(\{N+1, \dots, 2N\})} h(\omega \circ s).$$

Then we can write

$$1 - |1 - 2T_i(\sigma_i)| = (\sigma_i - \sigma_{i+N})^2 = \sigma_i + \sigma_{i+N} - 2\sigma_i\sigma_{i+N}.$$

Using this identity, we get for any exchangeable function $\lambda : \Omega \times \Theta \rightarrow \mathbb{R}_+$,

$$T \left\{ \exp \left[\lambda(\bar{r} - r_1) - \log \left[\cosh \left(\frac{\lambda}{2N} \right) \right] \sum_{i=1}^N (\sigma_i + \sigma_{i+N} - 2\sigma_i\sigma_{i+N}) \right] \right\} \leq 1.$$

Let us put

$$(3.2) \quad A(\lambda) = \frac{2N}{\lambda} \log \left[\cosh \left(\frac{\lambda}{2N} \right) \right],$$

$$(3.3) \quad v(\theta) = \frac{1}{2N} \sum_{i=1}^N (\sigma_i + \sigma_{i+N} - 2\sigma_i\sigma_{i+N}).$$

With this notation

$$T \left\{ \exp \left\{ \lambda [\bar{r} - r_1 - A(\lambda)v] \right\} \right\} \leq 1.$$

Let us notice now that

$$T'[v(\theta)] = \bar{r}(\theta) - r_1(\theta)r_2(\theta).$$

Let $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$ be any given exchangeable posterior distribution. Using the exchangeability of \mathbb{P} and π and the exchangeability of the exponential function, we get

$$\begin{aligned} & \mathbb{P} \left\{ \pi \left[\exp \left\{ \lambda [\bar{r} - r_1 - A(\bar{r} - r_1 r_2)] \right\} \right] \right\} = \mathbb{P} \left\{ \pi \left[\exp \left\{ \lambda [\bar{r} - r_1 - AT'(v)] \right\} \right] \right\} \\ & \leq \mathbb{P} \left\{ \pi \left[T' \exp \left\{ \lambda [\bar{r} - r_1 - Av] \right\} \right] \right\} = \mathbb{P} \left\{ T' \pi \left[\exp \left\{ \lambda [\bar{r} - r_1 - Av] \right\} \right] \right\} \\ & = \mathbb{P} \left\{ \pi \left[\exp \left\{ \lambda [\bar{r} - r_1 - Av] \right\} \right] \right\} = \mathbb{P} \left\{ T \pi \left[\exp \left\{ \lambda [\bar{r} - r_1 - Av] \right\} \right] \right\} \\ & = \mathbb{P} \left\{ \pi \left[T \exp \left\{ \lambda [\bar{r} - r_1 - Av] \right\} \right] \right\} \leq 1. \end{aligned}$$

We are thus ready to state

THEOREM 3.2.6. *In the case when $k = 1$, for any exchangeable probability distribution \mathbb{P} , for any exchangeable posterior distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$, for any exchangeable function $\lambda : \Omega \times \Theta \rightarrow \mathbb{R}_+$,*

$$\mathbb{P} \left\{ \exp \left[\sup_{\rho \in \mathcal{M}_+^1(\Theta)} \rho \left\{ \lambda [\bar{r} - r_1 - A(\lambda)(\bar{r} - r_1 r_2)] \right\} - \mathcal{K}(\rho, \pi) \right] \right\} \leq 1,$$

where $A(\lambda)$ is defined by equation (3.2, page 119).

We then deduce as previously

COROLLARY 3.2.7. *For any exchangeable posterior distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$, for any exchangeable probability measure $\mathbb{P} \in \mathcal{M}_+^1(\Omega)$, for any measurable exchangeable function $\lambda : \Omega \times \Theta \rightarrow \mathbb{R}_+$, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,*

$$\bar{r}(\theta) \leq r_1(\theta) + A(\lambda) [\bar{r}(\theta) - r_1(\theta) r_2(\theta)] - \frac{\log \{ \epsilon \pi [\Delta(\theta)] \}}{\lambda},$$

where $A(\lambda)$ is defined by equation (3.2, page 119).

In order to deduce an empirical bound from this theorem, we have to make some choice for $\lambda(\omega, \theta)$. Fortunately, it is easy to show that the bound holds uniformly in λ , because the inequality can be rewritten as a function of only one non-exchangeable quantity, namely $r_1(\theta)$. Indeed, since $r_2 = 2\bar{r} - r_1$, we see that the inequality can be written as

$$\bar{r}(\theta) \leq r_1(\theta) + A(\lambda) [\bar{r}(\theta) - 2\bar{r}(\theta) r_1(\theta) + r_1(\theta)^2] - \frac{\log \{ \epsilon \pi [\Delta(\theta)] \}}{\lambda}.$$

It can be solved in $r_1(\theta)$, to get

$$r_1(\theta) \geq f(\lambda, \bar{r}(\theta), -\log \{ \epsilon \pi [\Delta(\theta)] \}),$$

where

$$f(\lambda, \bar{r}, d) = [2A(\lambda)]^{-1} \left\{ 2\bar{r}A(\lambda) - 1 + \sqrt{[1 - 2\bar{r}A(\lambda)]^2 + 4A(\lambda) \left\{ \bar{r}[1 - A(\lambda)] - \frac{d}{\lambda} \right\}} \right\}.$$

Thus we can find some exchangeable function $\lambda(\omega, \theta)$, such that

$$f(\lambda(\omega, \theta), \bar{r}(\theta), -\log \{ \epsilon \pi [\Delta(\theta)] \}) = \sup_{\beta \in \mathbb{R}_+} f(\beta, \bar{r}(\theta), -\log \{ \epsilon \pi [\Delta(\theta)] \}).$$

Applying Corollary 3.2.7 (page 120) to that choice of λ , we see that

THEOREM 3.2.8. *For any exchangeable probability measure $\mathbb{P} \in \mathcal{M}_+^1(\Omega)$, for any exchangeable posterior probability distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$, for any $\lambda \in \mathbb{R}_+$,*

$$\bar{r}(\theta) \leq r_1(\theta) + A(\lambda) [\bar{r}(\theta) - r_1(\theta) r_2(\theta)] - \frac{\log \{ \epsilon \pi [\Delta(\theta)] \}}{\lambda},$$

where $A(\lambda)$ is defined by equation (3.2, page 119).

Solving the previous inequality in $r_2(\theta)$, we get

COROLLARY 3.2.9. *Under the same assumptions as in the previous theorem, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,*

$$r_2(\theta) \leq \inf_{\lambda \in \mathbb{R}_+} \frac{r_1(\theta) \left\{ 1 + \frac{2N}{\lambda} \log \left[\cosh \left(\frac{\lambda}{2N} \right) \right] \right\} - \frac{2 \log \{ \epsilon \pi [\Delta(\theta)] \}}{\lambda}}{1 - \frac{2N}{\lambda} \log \left[\cosh \left(\frac{\lambda}{2N} \right) \right] [1 - 2r_1(\theta)]}.$$

Applying this to our usual numerical example of a binary classification model with Vapnik–Cervonenkis dimension not greater than $h = 10$, when $N = 1000$, $\inf_{\Theta} r_1 = r_1(\hat{\theta}) = 10$ and $\epsilon = 0.01$, we obtain that $r_2(\hat{\theta}) \leq 0.4450$.

3.3. VAPNIK BOUNDS FOR INDUCTIVE CLASSIFICATION

3.3.1. ARBITRARY SHADOW SAMPLE SIZE. We assume in this section that

$$\mathbb{P} = \left(\bigotimes_{i=1}^N P_i \right)^{\otimes \infty} \in \mathcal{M}_+^1 \left\{ [(\mathcal{X} \times \mathcal{Y})^N]^{\mathbb{N}} \right\},$$

where $P_i \in \mathcal{M}_+^1(\mathcal{X} \times \mathcal{Y})$: we consider an infinite i.i.d. sequence of independent non-identically distributed samples of size N , the first one only being observed. More precisely, under \mathbb{P} each sample $(X_{i+jN}, Y_{i+jN})_{i=1}^N$ is distributed according to $\bigotimes_{i=1}^N P_i$, and they are all independent from each other. Only the first sample $(X_i, Y_i)_{i=1}^N$ is assumed to be observed. The shadow samples will only appear in the proofs. The aim of this section is to prove better Vapnik bounds, generalizing them in the same time to the independent non-i.i.d. setting, which to our knowledge has not been done before.

Let us introduce the notation $\mathbb{P}'[h(\omega)] = \mathbb{P}[h(\omega) | (X_i, Y_i)_{i=1}^N]$, where h may be any suitable (e.g. bounded) random variable, let us also put $\Omega = [(\mathcal{X} \times \mathcal{Y})^N]^{\mathbb{N}}$.

DEFINITION 3.3.1. For any subset $A \subset \mathbb{N}$ of integers, let $\mathfrak{C}(A)$ be the set of circular permutations of the totally ordered set A , extended to a permutation of \mathbb{N} by taking it to be the identity on the complement $\mathbb{N} \setminus A$ of A . We will say that a random function $h : \Omega \rightarrow \mathbb{R}$ is k -partially exchangeable if

$$h(\omega \circ s) = h(\omega), \quad s \in \mathfrak{C}(\{i + jN; j = 0, \dots, k\}), i = 1, \dots, N.$$

In the same way, we will say that a posterior distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$ is k -partially exchangeable if

$$\pi(\omega \circ s) = \pi(\omega) \in \mathcal{M}_+^1(\Theta), \quad s \in \mathfrak{C}(\{i + jN; j = 0, \dots, k\}), i = 1, \dots, N.$$

Note that \mathbb{P} itself is k -partially exchangeable for any k in the sense that for any bounded measurable function $h : \Omega \rightarrow \mathbb{R}$

$$\mathbb{P}[h(\omega \circ s)] = \mathbb{P}[h(\omega)], \quad s \in \mathfrak{C}(\{i + jN; j = 0, \dots, k\}), i = 1, \dots, N.$$

Let $\Delta_k(\theta) = \left\{ \theta' \in \Theta; [f_{\theta'}(X_i)]_{i=1}^{(k+1)N} = [f_{\theta}(X_i)]_{i=1}^{(k+1)N} \right\}$, $\theta \in \Theta, k \in \mathbb{N}^*$, and let

also $\bar{r}_k(\theta) = \frac{1}{(k+1)N} \sum_{i=1}^{(k+1)N} \mathbb{1}[f_{\theta}(X_i) \neq Y_i]$. Theorem 3.1.2 shows that for any

positive real parameter λ and any k -partially exchangeable posterior distribution $\pi_k : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$,

$$\mathbb{P} \left\{ \exp \left[\sup_{\theta \in \Theta} \lambda [\Phi_{\frac{\lambda}{N}}(\bar{r}_k) - r_1] + \log \{ \epsilon \pi_k [\Delta_k(\theta)] \} \right] \right\} \leq \epsilon.$$

Using the general fact that

$$\mathbb{P}[\exp(h)] = \mathbb{P} \left\{ \mathbb{P}'[\exp(h)] \right\} \geq \mathbb{P} \left\{ \exp[\mathbb{P}'(h)] \right\},$$

and the fact that the expectation of a supremum is larger than the supremum of an expectation, we see that with \mathbb{P} probability at most $1 - \epsilon$, for any $\theta \in \Theta$,

$$\mathbb{P}' \left\{ \Phi_{\frac{\lambda}{N}}[\bar{r}_k(\theta)] \right\} \leq r_1(\theta) - \frac{\mathbb{P}' \left\{ \log \{ \epsilon \pi_k [\Delta_k(\theta)] \} \right\}}{\lambda}.$$

For short let us put

$$\begin{aligned} \bar{d}_k(\theta) &= -\log \{ \epsilon \pi_k [\Delta_k(\theta)] \}, \\ d'_k(\theta) &= -\mathbb{P}' \left\{ \log \{ \epsilon \pi_k [\Delta_k(\theta)] \} \right\}, \\ d_k(\theta) &= -\mathbb{P} \left\{ \log \{ \epsilon \pi_k [\Delta_k(\theta)] \} \right\}. \end{aligned}$$

We can use the convexity of $\Phi_{\frac{\lambda}{N}}$ and the fact that $\mathbb{P}'(\bar{r}_k) = \frac{r_1 + kR}{k+1}$, to establish that

$$\mathbb{P}' \left\{ \Phi_{\frac{\lambda}{N}}[\bar{r}_k(\theta)] \right\} \geq \Phi_{\frac{\lambda}{N}} \left[\frac{r_1(\theta) + kR(\theta)}{k+1} \right].$$

We have proved

THEOREM 3.3.1. *Using the above hypotheses and notation, for any sequence $\pi_k : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$, where π_k is a k -partially exchangeable posterior distribution, for any positive real constant λ , any positive integer k , with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,*

$$\Phi_{\frac{\lambda}{N}} \left[\frac{r_1(\theta) + kR(\theta)}{k+1} \right] \leq r_1(\theta) + \frac{d'_k(\theta)}{\lambda}.$$

We can make as we did with Theorem 1.2.6 (page 11) the result of this theorem uniform in $\lambda \in \{\alpha^j; j \in \mathbb{N}^*\}$ and $k \in \mathbb{N}^*$ (considering on k the prior $\frac{1}{k(k+1)}$ and on j the prior $\frac{1}{j(j+1)}$), and obtain

THEOREM 3.3.2. *For any real parameter $\alpha > 1$, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,*

$$\begin{aligned} R(\theta) \leq & \frac{1 - \exp \left\{ -\frac{\alpha^j}{N} r_1(\theta) - \frac{1}{N} \left\{ d'_k(\theta) + \log [k(k+1)j(j+1)] \right\} \right\}}{\inf_{k \in \mathbb{N}^*, j \in \mathbb{N}^*} \frac{k}{k+1} \left[1 - \exp \left(-\frac{\alpha^j}{N} \right) \right]} \\ & - \frac{r_1(\theta)}{k}. \end{aligned}$$

As a special case we can choose π_k such that $\log\{\pi_k[\Delta_k(\theta)]\}$ is independent of θ and equal to $\log(\mathfrak{N}_k)$, where

$$\mathfrak{N}_k = \left| \left\{ [f_\theta(X_i)]_{i=1}^{(k+1)N}; \theta \in \Theta \right\} \right|$$

is the size of the trace of the classification model on the extended sample of size $(k+1)N$. With this choice, we obtain a bound involving a new flavour of conditional Vapnik entropy, namely

$$d'_k(\theta) = \mathbb{P}[\log(\mathfrak{N}_k) | (Z_i)_{i=1}^N] - \log(\epsilon).$$

In the case of binary classification using a Vapnik–Cervonenkis class of Vapnik–Cervonenkis dimension not greater than $h = 10$, when $N = 1000$, $\inf_{\Theta} r_1 = r_1(\hat{\theta}) = 0.2$ and $\epsilon = 0.01$, choosing $\alpha = 1.1$, we obtain $R(\hat{\theta}) \leq 0.4271$ (for an optimal value of $\lambda = 1071.8$, and an optimal value of $k = 16$).

3.3.2. A BETTER MINIMIZATION WITH RESPECT TO THE EXPONENTIAL PARAMETER. If we are not pleased with optimizing λ on a discrete subset of the real line, we can use a slightly different approach. From Theorem 3.1.2 (page 113), we see that for any positive integer k , for any k -partially exchangeable positive real measurable function $\lambda : \Omega \times \Theta \rightarrow \mathbb{R}_+$ satisfying equation (3.1, page 116) — with $\Delta(\theta)$ replaced with $\Delta_k(\theta)$ — for any $\epsilon \in (0, 1)$ and $\eta \in (0, 1)$,

$$\mathbb{P} \left\{ \mathbb{P}' \left[\exp \left[\sup_{\theta} \lambda \left[\Phi_{\frac{\lambda}{N}}(\bar{r}_k) - r_1 \right] + \log \{ \epsilon \eta \pi_k [\Delta_k(\theta)] \} \right] \right] \leq \epsilon \eta, \right.$$

therefore with \mathbb{P} probability at least $1 - \epsilon$,

$$\mathbb{P}' \left\{ \exp \left[\sup_{\theta} \lambda \left[\Phi_{\frac{\lambda}{N}}(\bar{r}_k) - r_1 \right] + \log \{ \epsilon \eta \pi_k [\Delta_k(\theta)] \} \right] \leq \eta, \right.$$

and consequently, with \mathbb{P} probability at least $1 - \epsilon$, with \mathbb{P}' probability at least $1 - \eta$, for any $\theta \in \Theta$,

$$\Phi_{\frac{\lambda}{N}}(\bar{r}_k) + \frac{\log \{ \epsilon \eta \pi_k [\Delta_k(\theta)] \}}{\lambda} \leq r_1.$$

Now we are entitled to choose

$$\lambda(\omega, \theta) \in \arg \max_{\lambda' \in \mathbb{R}_+} \Phi_{\frac{\lambda'}{N}}(\bar{r}_k) + \frac{\log \{ \epsilon \eta \pi_k [\Delta_k(\theta)] \}}{\lambda'}.$$

This shows that with \mathbb{P} probability at least $1 - \epsilon$, with \mathbb{P}' probability at least $1 - \eta$, for any $\theta \in \Theta$,

$$\sup_{\lambda \in \mathbb{R}_+} \Phi_{\frac{\lambda}{N}}(\bar{r}_k) - \frac{\bar{d}_k(\theta) - \log(\eta)}{\lambda} \leq r_1,$$

which can also be written

$$\Phi_{\frac{\lambda}{N}}(\bar{r}_k) - r_1 - \frac{\bar{d}_k(\theta)}{\lambda} \leq -\frac{\log(\eta)}{\lambda}, \quad \lambda \in \mathbb{R}_+.$$

Thus with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$, any $\lambda \in \mathbb{R}_+$,

$$\mathbb{P}' \left[\Phi_{\frac{\lambda}{N}}(\bar{r}_k) - r_1 - \frac{\bar{d}_k(\theta)}{\lambda} \leq -\frac{\log(\eta)}{\lambda} + \left[1 - r_1 + \frac{\log(\eta)}{\lambda} \right] \eta. \right.$$

On the other hand, $\Phi_{\frac{\lambda}{N}}$ being a convex function,

$$\begin{aligned} \mathbb{P}' \left[\Phi_{\frac{\lambda}{N}}(\bar{r}_k) - r_1 - \frac{\bar{d}_k(\theta)}{\lambda} \right] &\geq \Phi_{\frac{\lambda}{N}}[\mathbb{P}'(\bar{r}_k)] - r_1 - \frac{d'_k}{\lambda} \\ &= \Phi_{\frac{\lambda}{N}}\left(\frac{kR + r_1}{k + 1}\right) - r_1 - \frac{d'_k}{\lambda}. \end{aligned}$$

Thus with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,

$$\frac{kR + r_1}{k + 1} \leq \inf_{\lambda \in \mathbb{R}_+} \Phi_{\frac{\lambda}{N}}^{-1} \left[r_1(1 - \eta) + \eta + \frac{d'_k - \log(\eta)(1 - \eta)}{\lambda} \right].$$

We can generalize this approach by considering a finite decreasing sequence $\eta_0 = 1 > \eta_1 > \eta_2 > \dots > \eta_J > \eta_{J+1} = 0$, and the corresponding sequence of levels

$$\begin{aligned} L_j &= -\frac{\log(\eta_j)}{\lambda}, 0 \leq j \leq J, \\ L_{J+1} &= 1 - r_1 - \frac{\log(J) - \log(\epsilon)}{\lambda}. \end{aligned}$$

Taking a union bound in j , we see that with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$, for any $\lambda \in \mathbb{R}_+$,

$$\mathbb{P}' \left[\Phi_{\frac{\lambda}{N}}(\bar{r}_k) - r_1 - \frac{\bar{d}_k + \log(J)}{\lambda} \geq L_j \right] \leq \eta_j, \quad j = 0, \dots, J + 1,$$

and consequently

$$\begin{aligned} &\mathbb{P}' \left[\Phi_{\frac{\lambda}{N}}(\bar{r}_k) - r_1 - \frac{\bar{d}_k + \log(J)}{\lambda} \right] \\ &\leq \int_0^{L_{J+1}} \mathbb{P}' \left[\Phi_{\frac{\lambda}{N}}(\bar{r}_k) - r_1 - \frac{\bar{d}_k + \log(J)}{\lambda} \geq \alpha \right] d\alpha \leq \sum_{j=1}^{J+1} \eta_{j-1} (L_j - L_{j-1}) \\ &= \eta_J \left[1 - r_1 - \frac{\log(J) - \log(\epsilon) - \log(\eta_J)}{\lambda} \right] - \frac{\log(\eta_1)}{\lambda} + \sum_{j=1}^{J-1} \frac{\eta_j}{\lambda} \log \left(\frac{\eta_j}{\eta_{j+1}} \right). \end{aligned}$$

Let us put

$$\begin{aligned} d''_k[\theta, (\eta_j)_{j=1}^J] &= d'_k(\theta) + \log(J) - \log(\eta_1) \\ &\quad + \sum_{j=1}^{J-1} \eta_j \log \left(\frac{\eta_j}{\eta_{j+1}} \right) + \log \left(\frac{\epsilon \eta_J}{J} \right) \eta_J. \end{aligned}$$

We have proved that for any decreasing sequence $(\eta_j)_{j=1}^J$, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,

$$\frac{kR + r_1}{k + 1} \leq \inf_{\lambda \in \mathbb{R}_+} \Phi_{\frac{\lambda}{N}}^{-1} \left[r_1(1 - \eta_J) + \eta_J + \frac{d''_k[\theta, (\eta_j)_{j=1}^J]}{\lambda} \right].$$

REMARK 3.3.1. We can for instance choose $J = 2$, $\eta_2 = \frac{1}{10N}$, $\eta_1 = \frac{1}{\log(10N)}$, resulting in

$$d''_k = d'_k + \log(2) + \log \log(10N) + 1 - \frac{\log \log(10N)}{\log(10N)} - \frac{\log \left(\frac{20N}{\epsilon} \right)}{10N}.$$

In the case where $N = 1000$ and for any $\epsilon \in (0, 1)$, we get $d''_k \leq d'_k + 3.7$, in the case where $N = 10^6$, we get $d''_k \leq d'_k + 4.4$, and in the case $N = 10^9$, we get $d''_k \leq d'_k + 4.7$.

Therefore, for any practical purpose we could take $d''_k = d'_k + 4.7$ and $\eta_J = \frac{1}{10N}$ in the above inequality.

Taking moreover a weighted union bound in k , we get

THEOREM 3.3.3. *For any $\epsilon \in (0, 1)$, any sequence $1 > \eta_1 > \dots > \eta_J > 0$, any sequence $\pi_k : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$, where π_k is a k -partially exchangeable posterior distribution, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,*

$$R(\theta) \leq \inf_{k \in \mathbb{N}^*} \frac{k+1}{k} \inf_{\lambda \in \mathbb{R}_+} \Phi_{\frac{\lambda}{N}}^{-1} \left[r_1(\theta) + \eta_J [1 - r_1(\theta)] + \frac{d''_k[\theta, (\eta_j)_{j=1}^J] + \log[k(k+1)]}{\lambda} \right] - \frac{r_1(\theta)}{k}.$$

COROLLARY 3.3.4. *For any $\epsilon \in (0, 1)$, for any $N \leq 10^9$, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,*

$$R(\theta) \leq \inf_{k \in \mathbb{N}^*} \inf_{\lambda \in \mathbb{R}_+} \frac{k+1}{k} [1 - \exp(-\frac{\lambda}{N})]^{-1} \left\{ 1 - \exp \left[-\frac{\lambda}{N} \left[r_1(\theta) + \frac{1}{10N} \right] - \frac{\mathbb{P}'[\log(\mathfrak{N}_k) | (Z_i)_{i=1}^N] - \log(\epsilon) + \log[k(k+1)] + 4.7]}{N} \right] \right\} - \frac{r_1(\theta)}{k}.$$

Let us end this section with a numerical example: in the case of binary classification with a Vapnik–Cervonenkis class of dimension not greater than 10, when $N = 1000$, $\inf_{\Theta} r_1 = r_1(\hat{\theta}) = 0.2$ and $\epsilon = 0.01$, we get a bound $R(\hat{\theta}) \leq 0.4211$ (for optimal values of $k = 15$ and of $\lambda = 1010$).

3.3.3. EQUAL SHADOW AND TRAINING SAMPLE SIZES. In the case when $k = 1$, we can use Theorem 3.2.5 (page 118) and replace $\Phi_{\frac{\lambda}{N}}^{-1}(q)$ with $\left\{ 1 - \frac{2N}{\lambda} \times \log \left[\cosh \left(\frac{\lambda}{2N} \right) \right] \right\}^{-1} q$, resulting in

THEOREM 3.3.5. *For any $\epsilon \in (0, 1)$, any $N \leq 10^9$, any one-partially exchangeable posterior distribution $\pi_1 : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,*

$$R(\theta) \leq \inf_{\lambda \in \mathbb{R}_+} \frac{\left\{ 1 + \frac{2N}{\lambda} \log \left[\cosh \left(\frac{\lambda}{2N} \right) \right] \right\} r_1(\theta) + \frac{1}{5N} + 2 \frac{d'_1(\theta) + 4.7}{\lambda}}{1 - \frac{2N}{\lambda} \log \left[\cosh \left(\frac{\lambda}{2N} \right) \right]}.$$

3.3.4. IMPROVEMENT ON THE EQUAL SAMPLE SIZE BOUND IN THE I.I.D. CASE. Finally, in the case when \mathbb{P} is i.i.d., meaning that all the P_i are equal, we can improve the previous bound. For any partially exchangeable function $\lambda : \Omega \times \Theta \rightarrow \mathbb{R}_+$, we saw in the discussion preceding Theorem 3.2.6 (page 120) that

$$T \left[\exp[\lambda(\bar{r}_k - r_1) - A(\lambda)v] \right] \leq 1,$$

with the notation introduced therein. Thus for any partially exchangeable positive real measurable function $\lambda : \Omega \times \Theta \rightarrow \mathbb{R}_+$ satisfying equation (3.1, page 116), any one-partially exchangeable posterior distribution $\pi_1 : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$,

$$\mathbb{P} \left\{ \exp \left[\sup_{\theta \in \Theta} \lambda [\bar{r}_k(\theta) - r_1(\theta) - A(\lambda)v(\theta)] + \log[\epsilon \pi_1[\Delta(\theta)]] \right] \right\} \leq 1.$$

Therefore with \mathbb{P} probability at least $1 - \epsilon$, with \mathbb{P}' probability $1 - \eta$,

$$\bar{r}_k(\theta) \leq r_1(\theta) + A(\lambda)v(\theta) + \frac{1}{\lambda} [\bar{d}_1(\theta) - \log(\eta)].$$

We can then choose $\lambda(\omega, \theta) \in \arg \min_{\lambda' \in \mathbb{R}_+} A(\lambda')v(\theta) + \frac{\bar{d}_1(\theta) - \log(\eta)}{\lambda'}$, which satisfies the required conditions, to show that with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$, with \mathbb{P}' probability at least $1 - \eta$, for any $\lambda \in \mathbb{R}_+$,

$$\bar{r}_k(\theta) \leq r_1(\theta) + A(\lambda)v(\theta) + \frac{\bar{d}_1(\theta) - \log(\eta)}{\lambda}.$$

We can then take a union bound on a decreasing sequence of J values $\eta_1 \geq \dots \geq \eta_J$ of η . Weakening the order of quantifiers a little, we then obtain the following statement: with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$, for any $\lambda \in \mathbb{R}_+$, for any $j = 1, \dots, J$

$$\mathbb{P}' \left[\bar{r}_k(\theta) - r_1(\theta) - A(\lambda)v(\theta) - \frac{\bar{d}_1(\theta) + \log(J)}{\lambda} \geq -\frac{\log(\eta_j)}{\lambda} \right] \leq \eta_j.$$

Consequently for any $\lambda \in \mathbb{R}_+$,

$$\begin{aligned} \mathbb{P}' \left[\bar{r}_k(\theta) - r_1(\theta) - A(\lambda)v(\theta) - \frac{\bar{d}_1(\theta) + \log(J)}{\lambda} \right] \\ \leq -\frac{\log(\eta_1)}{\lambda} + \eta_J \left[1 - r_1(\theta) - \frac{\log(J) - \log(\epsilon) - \log(\eta_J)}{\lambda} \right] \\ + \sum_{j=1}^{J-1} \frac{\eta_j}{\lambda} \log \left(\frac{\eta_j}{\eta_{j+1}} \right). \end{aligned}$$

Moreover $\mathbb{P}'[v(\theta)] = \frac{r_1 + R}{2} - r_1 R$, (this is where we need equidistribution) thus proving that

$$\frac{R - r_1}{2} \leq \frac{A(\lambda)}{2} [R + r_1 - 2r_1 R] + \frac{d_1''[\theta, (\eta_j)_{j=1}^J]}{\lambda} + \eta_J [1 - r_1(\theta)].$$

Keeping track of quantifiers, we obtain

THEOREM 3.3.6. *For any decreasing sequence $(\eta_j)_{j=1}^J$, any $\epsilon \in (0, 1)$, any one-partially exchangeable posterior distribution $\pi : \Omega \rightarrow \mathcal{M}_+^1(\Theta)$, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,*

$$\begin{aligned} R(\theta) \leq \inf_{\lambda \in \mathbb{R}_+} \\ \frac{\left\{ 1 + \frac{2N}{\lambda} \log \left[\cosh \left(\frac{\lambda}{2N} \right) \right] \right\} r_1(\theta) + \frac{2d_1''[\theta, (\eta_j)_{j=1}^J]}{\lambda} + 2\eta_J [1 - r_1(\theta)]}{1 - \frac{2N}{\lambda} \log \left[\cosh \left(\frac{\lambda}{2N} \right) \right] [1 - 2r_1(\theta)]}. \end{aligned}$$

3.4. GAUSSIAN APPROXIMATION IN VAPNIK BOUNDS

3.4.1. GAUSSIAN UPPER BOUNDS OF VARIANCE TERMS. To obtain formulas which could be easily compared with original Vapnik bounds, we may replace $p - \Phi_a(p)$ with a Gaussian upper bound:

LEMMA 3.4.1. For any $p \in (0, \frac{1}{2})$, any $a \in \mathbb{R}_+$,

$$p - \Phi_a(p) \leq \frac{a}{2}p(1-p).$$

For any $p \in (\frac{1}{2}, 1)$,

$$p - \Phi_a(p) \leq \frac{a}{8}.$$

PROOF. Let us notice that for any $p \in (0, 1)$,

$$\begin{aligned} \frac{\partial}{\partial a} [-a\Phi_a(p)] &= -\frac{p \exp(-a)}{1-p+p \exp(-a)}, \\ \frac{\partial^2}{\partial^2 a} [-a\Phi_a(p)] &= \frac{p \exp(-a)}{1-p+p \exp(-a)} \left(1 - \frac{p \exp(-a)}{1-p+p \exp(-a)}\right) \\ &\leq \begin{cases} p(1-p) & p \in (0, \frac{1}{2}), \\ \frac{1}{4} & p \in (\frac{1}{2}, 1). \end{cases} \end{aligned}$$

Thus taking a Taylor expansion of order one with integral remainder:

$$-a\Phi(a) \leq \begin{cases} -ap + \int_0^a p(1-p)(a-b)db \\ \qquad \qquad \qquad = -ap + \frac{a^2}{2}p(1-p), & p \in (0, \frac{1}{2}), \\ -ap + \int_0^a \frac{1}{4}(a-b)db = -ap + \frac{a^2}{8}, & p \in (\frac{1}{2}, 1). \end{cases}$$

This ends the proof of our lemma. \square

LEMMA 3.4.2. Let us consider the bound

$$B(q, d) = \left(1 + \frac{2d}{N}\right)^{-1} \left[q + \frac{d}{N} + \sqrt{\frac{2dq(1-q)}{N} + \frac{d^2}{N^2}} \right], \quad q \in \mathbb{R}_+, d \in \mathbb{R}_+.$$

Let us also put

$$\bar{B}(q, d) = \begin{cases} B(q, d) & B(q, d) \leq \frac{1}{2}, \\ q + \sqrt{\frac{d}{2N}} & \text{otherwise.} \end{cases}$$

For any positive real parameters q and d

$$\inf_{\lambda \in \mathbb{R}_+} \Phi_{\frac{\lambda}{N}}^{-1} \left(q + \frac{d}{\lambda} \right) \leq \bar{B}(q, d).$$

PROOF. Let $p = \inf_{\lambda} \Phi_{\frac{\lambda}{N}}^{-1} \left(q + \frac{d}{\lambda} \right)$. For any $\lambda \in \mathbb{R}_+$,

$$p - \frac{\lambda}{2N} (p \wedge \frac{1}{2}) [1 - (p \wedge \frac{1}{2})] \leq \Phi_{\frac{\lambda}{N}}(p) \leq q + \frac{d}{\lambda}.$$

Thus

$$\begin{aligned} p &\leq q + \inf_{\lambda \in \mathbb{R}_+} \frac{\lambda}{2N} (p \wedge \tfrac{1}{2}) [1 - (p \wedge \tfrac{1}{2})] + \frac{d}{\lambda} \\ &= q + \sqrt{\frac{2d(p \wedge \tfrac{1}{2}) [1 - (p \wedge \tfrac{1}{2})]}{N}} \leq q + \sqrt{\frac{d}{2N}}. \end{aligned}$$

Then let us remark that $B(q, d) = \sup \left\{ p' \in \mathbb{R}_+; p' \leq q + \sqrt{\frac{2dp'(1-p')}{N}} \right\}$. If moreover $\frac{1}{2} \geq B(q, d)$, then according to this remark $\frac{1}{2} \geq q + \sqrt{\frac{d}{2N}} \geq p$. Therefore $p \leq \frac{1}{2}$, and consequently $p \leq q + \sqrt{\frac{2dp(1-p)}{N}}$, implying that $p \leq B(q, d)$. \square

3.4.2. ARBITRARY SHADOW SAMPLE SIZE. The previous lemma combined with Corollary 3.3.4 (page 125) implies

COROLLARY 3.4.3. *Let us use the notation introduced in Lemma 3.4.2 (page 127). For any $\epsilon \in (0, 1)$, any integer $N \leq 10^9$, with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,*

$$R(\theta) \leq \inf_{k \in \mathbb{N}^*} \frac{k+1}{k} \left\{ \bar{B} \left[r_1(\theta) + \frac{1}{10N}, d'_k(\theta) + \log[k(k+1)] + 4.7 \right] \right\} - \frac{r_1(\theta)}{k}.$$

3.4.3. EQUAL SAMPLE SIZES IN THE I.I.D. CASE. To make a link with Vapnik's result, it is useful to state the Gaussian approximation to Theorem 3.3.6 (page 126). Indeed, using the upper bound $A(\lambda) \leq \frac{\lambda}{4N}$, where $A(\lambda)$ is defined by equation (3.2) on page 119, we get with \mathbb{P} probability at least $1 - \epsilon$

$$R - r_1 - 2\eta_J \leq \inf_{\lambda \in \mathbb{R}_+} \frac{\lambda}{4N} [R + r_1 - 2r_1 R] + \frac{2d''_1}{\lambda} = \sqrt{\frac{2d''_1(R + r_1 - 2r_1 R)}{N}},$$

which can be solved in R to obtain

COROLLARY 3.4.4. *With \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$,*

$$\begin{aligned} R(\theta) &\leq r_1(\theta) + \frac{d''_1(\theta)}{N} [1 - 2r_1(\theta)] + 2\eta_J \\ &+ \sqrt{\frac{4d''_1(\theta) [1 - r_1(\theta)] r_1(\theta)}{N} + \frac{d''_1(\theta)^2}{N^2} [1 - 2r_1(\theta)]^2 + \frac{4d''_1(\theta)}{N} [1 - 2r_1(\theta)] \eta_J}. \end{aligned}$$

This is to be compared with Vapnik's result, as proved in Vapnik (1998, page 138):

THEOREM 3.4.5 (VAPNIK). *For any i.i.d. probability distribution \mathbb{P} , with \mathbb{P} probability at least $1 - \epsilon$, for any $\theta \in \Theta$, putting*

$$d_V = \log[\mathbb{P}(\mathfrak{N}_1)] + \log(4/\epsilon),$$

$$R(\theta) \leq r_1(\theta) + \frac{2d_V}{N} + \sqrt{\frac{4d_V r_1(\theta)}{N} + \frac{4d_V^2}{N^2}}.$$

Recalling that we can choose $(\eta_j)_{j=1}^2$ such that $\eta_J = \eta_2 = \frac{1}{10N}$ (which brings a negligible contribution to the bound) and such that for any $N \leq 10^9$,

$$d_1''(\theta) \leq \mathbb{P}[\log(\mathfrak{N}_1) | (Z_i)_{i=1}^N] - \log(\epsilon) + 4.7,$$

we see that our complexity term is somehow more satisfactory than Vapnik's, since it is integrated outside the logarithm, with a slightly larger additional constant (remember that $\log 4 \simeq 1.4$, which is better than our 4.7, which could presumably be improved by working out a better sequence η_j , but not down to $\log(4)$). Our variance term is better, since we get $r_1(1 - r_1)$, instead of r_1 . We also have $\frac{d_1''}{N}$ instead of $2\frac{d_V}{N}$, because we use no symmetrization trick.

Let us illustrate these bounds on a numerical example, corresponding to a situation where the sample is noisy or the classification model is weak. Let us assume that $N = 1000$, $\inf_{\Theta} r_1 = r_1(\hat{\theta}) = 0.2$, that we are performing binary classification with a model with Vapnik–Cervonenkis dimension not greater than $h = 10$, and that we work at confidence level $\epsilon = 0.01$. Vapnik's theorem provides an upper bound for $R(\hat{\theta})$ not smaller than 0.610, whereas Corollary 3.4.4 gives $R(\hat{\theta}) \leq 0.461$ (using the bound $d_1'' \leq d_1' + 3.7$ when $N = 1000$). Now if we go for Theorem 3.3.6 and do not make a Gaussian approximation, we get $R(\hat{\theta}) \leq 0.453$. It is interesting to remark that this bound is achieved for $\lambda = 1195 > N = 1000$. This explains why the Gaussian approximation in Vapnik's bound can be improved: for such a large value of λ , $\lambda r_1(\theta)$ does not behave like a Gaussian random variable.

Let us recall in conclusion that the best bound is provided by Theorem 3.3.3 (page 125), giving $R(\hat{\theta}) \leq 0.4211$, (that is approximately 2/3 of Vapnik's bound), for optimal values of $k = 15$, and of $\lambda = 1010$. This bound can be seen to take advantage of the fact that Bernoulli random variables are not Gaussian (its Gaussian approximation, Corollary 3.4.3, gives a bound $R(\theta) \simeq 0.4325$, still with an optimal $k = 15$), and of the fact that the optimal size of the shadow sample is significantly larger than the size of the observed sample. Moreover, Theorem 3.3.3 does not assume that the sample is i.i.d., but only that it is independent, thus generalizing Vapnik's bounds to inhomogeneous data (this will presumably be the case when data are collected from different places where the experimental conditions may not be the same, although they may reasonably be assumed to be independent).

Our little numerical example was chosen to illustrate the case when it is non-trivial to decide whether the chosen classifier does better than the 0.5 error rate of blind random classification. This case is of interest to choose “weak learners” to be aggregated or combined in some appropriate way in a second stage to reach a better classification rate. This stage of feature selection is unavoidable in many real world classification tasks. Our little computations are meant to exemplify the fact that Vapnik's bounds, although asymptotically suboptimal, as is obvious by comparison with the first two chapters, can do the job when dealing with moderate sample sizes.

