

A high-dimensional Wilks phenomenon

Stéphane Boucheron · Pascal Massart

February 9, 2010

Abstract A theorem by Wilks asserts that in smooth parametric density estimation the difference between the maximum likelihood and the likelihood of the sampling distribution converges toward a chi-square distribution where the number of degrees of freedom coincides with the model dimension. This observation is at the core of some goodness-of-fit testing procedures and of some classical model selection methods. This paper describes a non-asymptotic version of the Wilks phenomenon in bounded contrast optimization procedures. Using concentration inequalities for general functions of independent random variables, it proves that in bounded contrast minimization (as for example in Statistical Learning Theory), the difference between the empirical risk of the minimizer of the true risk in the model and the minimum of the empirical risk (the excess empirical risk) satisfies a Bernstein-like inequality where the variance term reflects the dimension of the model and the scale term reflects the noise conditions. From a mathematical statistics viewpoint, the significance of this result comes from the recent observation that when using model selection via penalization, the excess empirical risk represents a minimum penalty if non-asymptotic guarantees concerning prediction error are to be provided. From the perspective of empirical process theory, this paper describes a concentration inequality for the supremum of a bounded non-centered (actually non-positive) empirical process. Combining the now classical analysis of M-estimation (building on Talagrand's inequality for suprema of empirical processes) and versatile moment inequalities for functions of independent random variables, this paper develops a genuine Bernstein-like inequality that seems beyond the reach of traditional tools.

Supported by ANR grant TAMIS

Supported by Network of Excellence PASCAL II

Stéphane Boucheron
Laboratoire Probabilités et Modèles Aléatoires
Université Paris-Diderot
175 rue du Chevaleret, 75013-F Paris
E-mail: stephane.boucheron@math.jussieu.fr

Pascal Massart
Département de Mathématiques
Université Paris-Sud
91405-F Orsay
E-mail: pascal.massart@math.u-psud.fr

Keywords Wilks phenomenon · Risk estimates · Suprema of empirical processes · Concentration inequalities · Statistical learning

Mathematics Subject Classification (2000)

1 Introduction

Soon after Fisher popularized maximum likelihood methods, Wilks described the so-called Wilks phenomenon (Wilks 1938). Assume $(P_\theta, \theta \in \Theta \subseteq \mathbb{R}^m)$ defines a family of probability distributions over some space \mathcal{X} where for each θ , p_θ is a version of the density of P_θ with respect to some dominating measure μ . Given a sample x_1, \dots, x_n , of elements of \mathcal{X} , the log-likelihood at θ is defined as $\ell_n(\theta) = \sum_{i=1}^n \log p_\theta(x_i)$. Assume that $\mu^{\otimes n}$ -almost surely, there exists a maximizer $\hat{\theta} \in \Theta$ of the log-likelihood, that is

$$\ell_n(\hat{\theta}) = \sup_{\theta \in \Theta} \sum_{i=1}^n \log p_\theta(x_i).$$

Then, if the model is smooth enough (see Bickel and Doksum 1976; van der Vaart 1998, for precise definitions and demonstrations), assuming that X_1, \dots, X_n, \dots are sampled independently according to P_θ , where $\theta \in \Theta$, the sequences of random variables $2(\ell_n(\hat{\theta}) - \ell_n(\theta))$ and $2nD(P_\theta, P_{\hat{\theta}})$ both converge in distribution toward χ_m^2 (chi-square distribution with m degrees of freedom) where $D(P, Q)$ denotes the relative entropy or Kullback information between distributions P and Q . This remarkable phenomenon can be used to build goodness-of-fit tests. In Gaussian linear models with known variance, the Wilks phenomenon is used to perform variable selection (although in disguised form) (Akaike 1974; Mallows 1973).

Considering maximum likelihood estimation as a special case of M-estimation, one may wonder whether extensions of the Wilks phenomenon exist. Does it hold in non-parametric density estimation problems where maximum likelihood estimation or generalizations of it can still be used? Does it still hold when dealing with other contrast optimization problems? A short glimpse at the proof of the Wilks phenomenon reveals that it depends on the asymptotic normality of $\sqrt{n}(\hat{\theta} - \theta)$. When this last property does not hold, for example in non-parametric inference, we should look for reasonable generalizations. Such an endeavor is the main topic of the paper entitled ‘‘Generalized likelihood ratio statistics and Wilks phenomenon’’ by Fan et al. (2001). Though our aims are quite different, the findings of Fan et al. help to motivate our approach. They start by exploring a simple non-parametric Gaussian regression model where the parameter space is a Sobolev ball. The authors aim at testing whether the regression function is affine or not. They analyze a statistics that compares the maximum log-likelihood in the Sobolev ball and the maximum log-likelihood in the submodel constituted by affine functions, and generalizations of this approach. The traditional Wilks phenomenon does not hold in such an infinite-dimensional setting. But, Fan et al. point out that the number of non-zero coefficients of the maximum likelihood estimator (the effective dimension of the maximum-likelihood estimator) inside the Sobolev ball tends to increase with sample size. Meanwhile, as m tends to infinity, if $X_m \sim \chi_m^2$ ($(X_m - \mathbb{E}[X_m])/\sqrt{\mathbb{E}[X_m]}$ converges in distribution towards a centered Gaussian random variable with variance 2). This observation allows Fan et al. to propose a convenient generalization of the Wilks phenomenon: when centered, the difference between the maximum log-likelihood and the likelihood in the submodel should converge toward a Gamma distribution whose shape parameter depends mostly on the effective dimension of the maximum likelihood estimator and whose scale parameter is slowly varying with respect to sample size

and effective dimension. This asymptotic pivotality property paves the way to non-trivial statistical applications. A generalized Wilks phenomenon has been reported in several settings (Fan 1993; Portnoy 1988; Gayraud and Pouet 2001). But Fan et al. observe:

While we have observed the Wilks phenomenon and demonstrated it for a few useful cases, it is impossible to verify the phenomenon for all non-parametric hypothesis problems. The Wilks phenomenon needs to be checked for other problems not covered in (Fan et al. 2001).

Rather than attempting to prove asymptotic distributional results, we are looking for non-asymptotic moment or tail bounds concerning analogues of $\ell_n(\hat{\theta}) - \ell_n(\theta)$ that exhibit a Gamma-like behavior where the shape parameter reflects the model dimension rather than the sample size and the scale parameter depends only moderately on model dimension and sample size. To get a more concrete understanding, remember that if Z is a $\Gamma(p, 1/c)$ -distributed random variable with expectation $\mathbb{E}[Z] = cp$ and variance $\sigma^2 = c^2p$, then for $q \geq 2$, the q -th norm of $(Z - \mathbb{E}[Z])_+$ can be upper-bounded by $\kappa(\sqrt{q\sigma^2} + cq)$ where κ is a universal constant that depends neither on the expectation nor on the variance of the Γ distribution.

The settings considered in (Portnoy 1988; Fan 1993; Fan et al. 2001) and the references we are aware of share a common feature: they are connected with density estimation in a Gaussian framework. They disregard robustness considerations. This contrasts with the setting we are interested in: statistical learning (Vapnik. 1995; Schoelkopf and Smola 2002). In statistical learning, a product set $\mathcal{X} \times \mathcal{Y}$ is endowed with a probability distribution P , the components of $\omega = (x, y) \in \mathcal{X} \times \mathcal{Y}$ are denoted by X ($X(\omega) = x$) and Y ($Y(\omega) = y$). In binary classification, $\mathcal{Y} = \{-1, 1\}$. A loss function ℓ maps $\mathcal{Y} \times \mathcal{Y}$ to \mathbb{R}_+ . In the sequel, all loss functions are assumed to be bounded. For example, the classification loss is defined by $\ell(y, y') = \mathbf{1}_{y \neq y'}$ (Vapnik. 1995).

In the sequel, if U is distributed according to P , then for any (measurable) function g , the expected value of $g(U)$ is denoted either by Pg , $Pg(U)$, $\mathbb{E}_P[g(U)]$ or even $\mathbb{E}[g(U)]$ when P is clear from context. If U_1, \dots, U_n is a collection of independent random variables, P_n denotes the associated empirical distribution: $P_n g = \frac{1}{n} \sum_{i=1}^n g(U_i)$.

The learning problem consists in finding a function f on \mathcal{X} such that the risk $R(f) = P\ell(f(X), Y)$ is as small as possible. Henceforth f^* is assumed to minimize $R(f)$ among all functions such that $(x, y) \mapsto \ell(f(x), y)$ is measurable. For example, in binary classification, the best possible classifier f^* which is called the Bayes classifier is defined from the regression function $\eta(x) = \mathbb{E}[Y | X = x]$ by $f^*(x) = \text{sign}(\eta(x))$.

A learning algorithm typically looks for a good approximation \hat{f} of f^* in some class \mathcal{F} of classifiers by minimizing an empirical risk like $R_n(f) = P_n\ell(f(X), Y)$. We do not assume that $f^* \in \mathcal{F}$ but we assume that some $\bar{f} \in \mathcal{F}$ minimizes the risk $R(\cdot)$ over \mathcal{F} . Henceforth, the model bias is $R(\bar{f}) - R(f^*)$. The connections between the richness of \mathcal{F} , the model bias, the distribution of $|\eta(\cdot) - 1/2|$ (the so-called noise conditions from Tsybakov (2004)), and the distribution of the risk $R(\hat{f}) = P\ell(\hat{f}(X), Y)$ and more recently the excess risk $R(\hat{f}) - R(\bar{f})$ have been the subject of intense research during the last thirty-five years (Vapnik. 1995; Mammen and Tsybakov 1999; Massart 2000b; Koltchinskii 2006; Bartlett and Mendelson 2006; Massart 2007; Massart and Nédélec 2006). The present paper is concerned with a rather atypical object, the empirical analogue of the excess risk, that from now on, will be called the unnormalized *empirical excess risk*:

$$n(R_n(\bar{f}) - R_n(\hat{f})) = n \sup_{f \in \mathcal{F}} (R_n(\bar{f}) - R_n(f)).$$

The purpose of the paper is to show that the right-tail bounds for the excess empirical risk described in (Koltchinskii 2006; Massart 2007) can be refined into concentration inequalities.

The main results of the paper are stated in Sections 3 and 4. Upper-bounds on the variance of the empirical excess risk are first derived by building on the Efron-Stein inequalities (Inequality (2)) and classical tail bounds for excess risk and excess empirical risk (Theorem 1). Those variance estimates (Theorem 2) already show the merits of an ad hoc analysis of this atypical supremum of an empirical process. Variance estimates are illustrated in a simple setting: bounded regression using piecewise constant functions (histograms) with quadratic loss. It has been possible to check that the expected empirical excess risk is almost a linear function of the model dimension (that is the number D of pieces of the partition that defines the model) (Arlot and Massart 2009): under mild conditions on the relationship between the number of classes and the sample size, $\mathbb{E}Z/\sigma^2 \approx D$ where $Z = n(R_n(\bar{f}) - R_n(\hat{f}_n))$. And in this simple context, Theorem 2 implies that the variance of the unnormalized empirical excess risk is upper-bounded by $\kappa\sigma^2 D$ where κ is some universal constant.

These variance estimates are amplified In Section 4 in order to estimate the higher moments of the centered empirical excess risk. Taking advantage on moment estimates for general functions of independent random variables derived in (Boucheron et al. 2005b), it is possible to show that the unnormalized empirical excess risk satisfies a Bernstein-like moment inequality where the variance term is of the order of the model dimension, while the scale term does not depend too much on either sample size or model dimension.

The end of Section 4 turns back to the regression problem. It follows that for a fixed noise level σ , there exists a universal constant κ such that the unnormalized empirical excess risk satisfies:

$$\|Z - \mathbb{E}Z\|_q \leq \kappa(\sqrt{Dq} + q)$$

for all sample sizes n and all partition cardinalities D where κ is some universal constant. In this simplified setting, the empirical excess risk exemplifies the high dimensional Wilks phenomenon. The significance of this observation is best understood by having a second look at the already mentioned paper by Arlot and Massart (2009) on model selection. The concentration of the empirical excess risk is the keystone in the derivation of an adaptive model selection method: it justifies a data-dependent penalty calibration technique.

Section 5 describes what can be obtained from theorems from Section 4 in classification problems. When using Vapnik-Chervonenkis classes of classifiers, the VC-dimension is a well-known surrogate for the traditional metric dimension that appears in the analysis of smooth parametric models. Under the so-called random classification noise (Angluin and Laird 1987; Kearns et al. 1997; Bartlett et al. 2002), noise conditions turn out to be very simple and the sharp estimates on the excess risk derived by Massart and Nédélec (2006) can be injected in Theorems 2 and 5. In the absence of model bias, variance and higher moment upper-bounds for excess empirical risk provide another version of the Wilks phenomenon: variance is upper-bounded by the VC-dimension multiplied by a slowly varying function of the ratio between sample size and VC-dimension. And, as far as the higher moment estimates are concerned, the scale parameter only depends on the noise conditions.

2 Notation and background

The expression $n(R_n(\bar{f}) - R_n(\hat{f})) = n \sup_{f \in \mathcal{F}} (R_n(\bar{f}) - R_n(f))$ stresses the fact that the empirical excess risk is the supremum of an empirical process. The latter is atypical

since the expected value of the supremum is always non-negative while for each $f \in \mathcal{F}$, $\mathbb{E} [R_n(\bar{f}) - R_n(f)] \leq 0$. We are not aware of any direct method that could be used to analyze the excess empirical risk.

In the sequel, tools from empirical process theory are not directly applied to \mathcal{F} but rather to the associated (relative) loss class \mathcal{H} . The loss class \mathcal{H} generated by \mathcal{F} is defined as

$$\mathcal{H} = \{(x, y) \mapsto \ell(\bar{f}(x), y) - \ell(f(x), y) : f \in \mathcal{F}\}.$$

Henceforth h^* denotes the loss associated with the global risk minimizer f^*

$$h^*(x, y) = \ell(\bar{f}(x), y) - \ell(f^*(x), y),$$

and \hat{h}_n denotes the loss function associated with the empirical risk minimizer \hat{f} :

$$\hat{h}_n(x, y) = \ell(\bar{f}(x), y) - \ell(\hat{f}(x), y).$$

When n is clear from context, we may abbreviate \hat{h}_n by \hat{h} . The excess empirical risk is thus rewritten as a supremum of an empirical process: $n(R_n(\bar{f}) - R_n(\hat{f})) = nP_n\hat{h}_n = n \sup_{h \in \mathcal{H}} P_n h$. Note that the model bias $R(\bar{f}) - R(f^*)$ is equal to $P h^*$ and that the excess risk $R(\hat{f}) - R(\bar{f})$ can be rewritten as $-P\hat{h}$.

For the sake of simplicity, in the sequel, we assume that \mathcal{H} is separable, that is \mathcal{H} contains a dense countable subset \mathcal{G} (with respect to the topology of pointwise convergence). Since \mathcal{H} is also assumed to be uniformly bounded, dominated convergence arguments ensure for example that:

$$\sup_{h \in \mathcal{H}} (P_n - P)h = \sup_{h \in \mathcal{G}} (P_n - P)h.$$

This ensures that the supremum over \mathcal{H} is measurable. Suprema over intersection of \mathcal{H} with various balls are checked to be measurable by similar arguments. We assume furthermore that all suprema over \mathcal{H} are achieved. Dealing with approximate minimizers of the empirical risk makes proof slightly lengthier without providing more insights.

We need to recall a few important notions that allow us to derive tail-bounds for excess risk and empirical excess risk. Those tail bounds are not meant to precisely describe the fluctuations of the empirical excess risk around its expectations. This is the purpose of the concentration inequalities stated in this paper.

Thanks to ideas borrowed from robust statistics (Huber 1967), empirical process theory (Shorack and Wellner 1986; van der Vaart and Wellner 1996; van de Geer 2000) and the theory of concentration inequalities (Talagrand 1996a,b; Ledoux 2001; Massart 2007), we now have a rather precise idea of the aforementioned connexions (Boucheron et al. 2005a; Massart 2007, for surveys).

In the sequel, various properties of the following class of positive sub-linear functions are repeatedly used. The class \mathcal{C}_1 comprises non-decreasing and continuous functions ϕ from \mathbb{R}_+ to \mathbb{R}_+ such that $\phi(x)/x$ is non-increasing and $\phi(1) \geq 1$. Note that if ψ and ρ belong to \mathcal{C}_1 , so do $\psi \circ \rho$ and $\psi + \rho$. Moreover, for any $\alpha \geq \psi(1)$, the equation $\alpha r^2 = \psi(r)$ has a unique solution in $(0, 1]$. One can check that any function $\psi \in \mathcal{C}_1$, is sub-additive: $\psi(u+v) \leq \psi(u) + \psi(v)$. Note also that if U is distributed according to P , $\mathbb{E}_P[\psi(U)] \leq 2\psi(\mathbb{E}_P[U])$.

In the sequel, we assume that there exists a pseudo-distance d on the set of bounded real functions on \mathcal{X} , such that for $f \in \mathcal{F}$

$$\mathbb{E}_P \left[(\ell(f(X), Y) - \ell(f^*(X), Y))^2 \right] \leq d^2(f, f^*), \quad (1)$$

or equivalently, using relative loss functions: $P(h^* - h)^2 \leq d^2(f, f^*)$. In some situations, as when analyzing classification problems, it is possible to define d as $d^2(f, f^*) = P(h^* - h)^2$, but in other circumstances, for example when investigating bounded regression, it is convenient to define d in other ways.

The complexity assumption below aims at describing the richness of the L_2 neighborhood of 0 in the loss class \mathcal{H} .

Assumption 1 [COMPLEXITY ASSUMPTION] *There exists some function $\psi \in \mathcal{C}_1$ such that for all r and all integers n satisfying $\psi(r) \leq \sqrt{nr^2}$:*

$$\sqrt{n}\mathbb{E} \left[\sup_{f \in \mathcal{F}: d^2(f, \bar{f}) \leq r^2} |(P_n - P)(\ell(f(X), Y) - \ell(\bar{f}(X), Y))| \right] \leq \psi(r).$$

The possibility to derive interesting bounds (fast rates) for the excess risk critically relies on the next assumption (Mammen and Tsybakov 1999; Massart 2000b; Koltchinskii 2006; Bartlett and Mendelson 2006; Massart 2007; Boucheron et al. 2005a). The name “noise conditions” will be clarified through examples in Sections 3, 4 and 5.

Assumption 2 [NOISE CONDITIONS] *There exists some function $\omega \in \mathcal{C}_1$ such that for all f in \mathcal{F} (with associated relative loss function $h \in \mathcal{H}$):*

$$d^2(f, f^*) \leq \omega^2 \left(\sqrt{P(h^* - h)} \right) = \omega^2 \left(\sqrt{R(f) - R(f^*)} \right).$$

Throughout the paper, a *learning task* is defined by a sampling probability P over $\mathcal{X} \times \mathcal{Y}$, a loss function ℓ , and a collection of functions \mathcal{F} (a relative loss class \mathcal{H} is automatically associated with the learning task). If the complexity and noise conditions assumptions are satisfied, the positive root r_* of the equation

$$\sqrt{nr^2} = \psi(\omega(r))$$

is a key quantity (van der Vaart and Wellner 1996; van de Geer 2000; van der Vaart 1998). It appears in upper-bounds for the expected excess risk $\mathbb{E}[-P\hat{h}_n]$ and the expected empirical risk $\mathbb{E}[P_n\hat{h}_n]$. Moreover, the analysis reveals that a function of r_* upper-bounds the expected value of $P\hat{h}_n^2$ and $P_n\hat{h}_n^2$ (Theorem 1 and Theorem 4).

The positive root r_* of equation $\sqrt{nr^2} = \psi(\omega(r))$ also shows up in the detailed analysis of our main object of concern, the unnormalized empirical excess risk $nP_n\hat{h}_n$.

3 The variance of the empirical excess risk

3.1 A generic upper-bound

The first step in understanding the fluctuations of the empirical excess risk around its expectation consists of upper-bounding the variance and checking that it can be related with the expectation. In order to carry out this program we use the Efron-Stein inequalities.

Let Z denote a (square-integrable) function of a sequence of independent random variables (U_1, U_2, \dots, U_n) , that is $Z = F(U_1, \dots, U_n)$ for some function F . Let U'_1, \dots, U'_n be distributed as U_1, \dots, U_n and be independent of U_1, \dots, U_n . For each $i \leq n$, let $Z'_i = F(U_1, \dots, U_{i-1}, U'_i, U_{i+1}, \dots, U_n)$, let $U^{(i)} = (U_1, \dots, U_{i-1}, U_{i+1}, \dots, U_n)$. Let F_i denote

a function of $n - 1$ arguments and let $Z_i = F_i(U_1, \dots, U_{i-1}, U_{i+1}, U_n) = F_i(U^{(i)})$. The jackknife estimates of variance V_+ and V (Quenouille 1949; Tukey 1958) are defined by

$$V_+ = \sum_{i=1}^n \mathbb{E} \left[(Z - Z'_i)_+^2 \mid U_1, \dots, U_n \right] \quad \text{and} \quad V = \sum_{i=1}^n (Z - Z_i)^2 .$$

The Efron-Stein inequalities (Efron and Stein 1981) assert that the expectation of the jackknife estimates of variance are upper-bounds:

$$\text{Var}[Z] \leq \mathbb{E}[V_+] \leq \mathbb{E}[V] . \quad (2)$$

Before handling the variance of the empirical excess risk, let us recall the following upper-bounds on the variance of suprema of general empirical processes. In the statement of the next proposition \mathcal{H} denotes a generic class of bounded functions over some space \mathcal{X} . The random variables U_i are assumed to be independent and \mathcal{X} -valued. The following inequality follows from $\text{Var}[Z] \leq \mathbb{E}[V_+]$: let $Z = n \sup_{h \in \mathcal{H}} P_n h$, then

$$\text{Var}[Z] \leq 2\mathbb{E} \left[n \sup_{h \in \mathcal{H}} P_n h^2 \right] + 2 \sup_{h \in \mathcal{H}} \mathbb{E} \left[\sum_{i=1}^n h^2(U_i) \right] . \quad (3)$$

The upper-bound can be further simplified and replaced by $4\mathbb{E} [n \sup_{h \in \mathcal{H}} P_n h^2]$.

When the random variables $h(U_i)$ are centered and identically distributed, the following handy bound can be derived as in (Rio 2001):

$$\text{Var}[Z] \leq 2b\mathbb{E}[Z] + n\sigma^2 \quad (4)$$

where $b = \sup_{h \in \mathcal{H}} \sup_x |h(x)|$ and $\sigma^2 = \sup_{h \in \mathcal{H}} \text{Var} [h(U)]$.

As the empirical excess risk is also the supremum of an empirical process, it is not too surprising that the Efron-Stein inequalities allow us to relate the variance of the empirical excess risk with its expectation and with various functions of the expected excess risk.

Proposition 1 *Let \mathcal{H} denote a class of measurable functions over \mathcal{U} .*

Let P denote a probability distribution over \mathcal{U} . Assume that functions in \mathcal{H} are square-integrable. Let \hat{h}_n maximize $P_n h$ over \mathcal{H} . Then

$$\frac{\text{Var} [n P_n \hat{h}_n]}{2n} \leq \mathbb{E} [P_n \hat{h}_n^2] + \mathbb{E} [P \hat{h}_n^2] . \quad (5)$$

Assume now that for all $h \in \mathcal{H}$, $P h \leq 0$, that the null function belongs to \mathcal{H} and that $\|h\|_\infty \leq 1$, then

$$\frac{\text{Var} [n P_n \hat{h}_n]}{n} \leq 2\mathbb{E} [(P_{n-1} - P) \hat{h}_{n-1}] + \mathbb{E} [P \hat{h}_{n-1}^2] . \quad (6)$$

Proof Let (U_1, \dots, U_n) be independently, identically distributed according to P . Let (U'_1, \dots, U'_n) denote an independent copy of (U_1, \dots, U_n) . The first inequality in the proposition is obtained by using the first Efron-Stein inequality ($\text{Var}[Z] \leq \mathbb{E}[V_+]$), and taking advantage of the fact that for each $i \in 1, \dots, n$

$$(Z - Z'_i) \leq \hat{h}_n(U_i) - \hat{h}_n(U'_i) .$$

Thus, as $(a + b)^2 \leq 2(a^2 + b^2)$, taking expectations over primed variables:

$$V_+ \leq 2n \left(P_n \widehat{h}_n^2 + P \widehat{h}_n^2 \right).$$

Taking expectations over U_1, \dots, U_n :

$$\text{Var} \left[n P_n \widehat{h}_n \right] \leq \mathbb{E}[V_+] \leq 2n \mathbb{E} \left[P_n \widehat{h}_n^2 \right] + 2n \mathbb{E} \left[P \widehat{h}_n^2 \right].$$

Let us now turn to the proof of the second inequality. Let $\widehat{h}_n^{(i)}$ denote the maximizer of the empirical process when the i^{th} element has been removed from the sample, that is $\widehat{h}_n^{(i)}$ maximizes $\sum_{j \leq n, j \neq i} h(U_j)$. Thus

$$Z = \sum_{i \leq n} \widehat{h}_n(U_i) \quad \text{and let} \quad Z_i = \sum_{j \leq n, j \neq i} \widehat{h}_n^{(i)}(U_j).$$

Note first that for each $i \in 1, \dots, n$

$$-1 \leq \widehat{h}_n^{(i)}(U_i) \leq Z - Z_i \leq \widehat{h}_n(U_i) \leq 1,$$

hence, the so-called self-boundedness property holds: $\sum_{i=1}^n (Z - Z_i) \leq Z$. For each sample (u_1, \dots, u_n) , for each i ,

$$\begin{aligned} (Z - Z_i)^2 - \left(\widehat{h}_n^{(i)}(u_i) \right)^2 &= \left(Z - Z_i + \widehat{h}_n^{(i)}(u_i) \right) \left(Z - Z_i - \widehat{h}_n^{(i)}(u_i) \right) \\ &\quad \text{and since the second factor is non-negative} \\ &\leq 2 \left(Z - Z_i - \widehat{h}_n^{(i)}(u_i) \right). \end{aligned}$$

Summing over all i , using the self-boundedness property:

$$\sum_{i=1}^n (Z - Z_i)^2 \leq 2Z + \sum_{i=1}^n \left(\widehat{h}_n^{(i)}(u_i) \right)^2 - 2 \sum_{i=1}^n \widehat{h}_n^{(i)}(u_i).$$

Now, taking expectations and invoking the Efron-Stein inequality:

$$\begin{aligned} \text{Var}[Z] &\leq 2\mathbb{E}[Z] + \sum_{i=1}^n \mathbb{E} \left[\left(\widehat{h}_n^{(i)}(U_i) \right)^2 \right] - 2 \sum_{i=1}^n \mathbb{E} \left[\widehat{h}_n^{(i)}(U_i) \right] \\ &= 2\mathbb{E}[Z] + n \mathbb{E} \left[\left(\widehat{h}_n^{(n)}(U_n) \right)^2 \right] - 2n \mathbb{E} \left[P \widehat{h}_n^{(n)} \right] \quad \text{by exchangeability} \\ &= 2n \mathbb{E}[P_n \widehat{h}_n] + n \mathbb{E} \left[P \widehat{h}_{n-1}^2 \right] - 2n \mathbb{E} \left[P \widehat{h}_{n-1} \right]. \end{aligned}$$

It remains to check that $\mathbb{E}[P_n \widehat{h}_n] \leq \mathbb{E}[P_{n-1} \widehat{h}_{n-1}]$. This follows from a classical argument that works for suprema of empirical processes provided some measurability conditions are met. Let \mathcal{G}_n denote the sub- σ -algebra of $\sigma(U_i, i \in \mathbb{N})$ consisting of events that are invariant under permutation of the first n coordinates. Then P_n is a measure-valued backward martingale: for any $f \in \mathcal{H}$, $\mathbb{E}[P_{n-1} h \mid \mathcal{G}_n] = P_n h$ a.s. As a supremum of linear functionals,

$P_n \widehat{h}_n$ is a backward submartingale with respect to the decreasing filtration (\mathcal{G}_n) (Pollard 1984). Almost surely, we have

$$\begin{aligned} \mathbb{E} \left[P_{n-1} \widehat{h}_{n-1} \mid \mathcal{G}_n \right] &= \mathbb{E} \left[\sup_{h \in \mathcal{H}} P_{n-1} h \mid \mathcal{G}_n \right] \\ &\geq \sup_{h \in \mathcal{H}} \mathbb{E} [P_{n-1} h \mid \mathcal{G}_n] \\ &= \sup_{h \in \mathcal{H}} P_n h \\ &= P_n \widehat{h}_n . \end{aligned}$$

So that by deconditioning $\mathbb{E} \left[P_{n-1} \widehat{h}_{n-1} \right] \geq \mathbb{E} \left[P_n \widehat{h}_n \right]$. \square

Remark 1 Note that we refrained from using the immediate upper-bounds

$$P_n \widehat{h}_n^2 \leq \sup_{h \in \mathcal{H}} P_n h^2 \quad \text{or} \quad P \widehat{h}_n^2 \leq \sup_{h \in \mathcal{H}} P h^2$$

as is usually done in the proof of the above-stated upper-bounds on the variance of suprema of general empirical processes (Rio 2001; Bousquet 2003). This would lead to bounds like (3) and would not allow us to take advantage of the fact that, under favorable noise conditions, with high probability, the supremum of the empirical process is achieved at some \widehat{f}_n in the L_2 neighborhood of \bar{f} .

3.2 Learning rates and variance upper-bounds

Up to now, it is not clear whether cautiousness is to be rewarded. In order to check that this is the case, we turn back to statistical learning problems, and invoke Proposition 1 when the function class denoted by \mathcal{H} in Proposition 1 turns out to be the loss class associated with a learning task as explained on page 6. In order to take full advantage of Proposition 1, we need non-trivial bounds on the expected values of what becomes the excess risk $-P \widehat{h}_n$ and the excess empirical risk $P_n \widehat{h}_n$.

Remark 2 When \mathcal{H} is the loss class associated with a learning task, inequality (5) stresses the fact that the variance of the empirical excess risk is upper-bounded using the second moments of the $L_2(P)$ and the $L_2(P_n)$ distances between the risk minimizer and the empirical risk minimizer. A similar connection between the so-called L_1 -stability $\mathbb{E} \left[P |\widehat{h}_n - \widehat{h}_{n-1}| \right]$ and the bias of the leave-one-out estimate of the risk had been pointed out by Devroye and Wagner (1977) and further investigated by Kearns and Ron (1999) and Rakhlin et al. (2005) (see Theorem 4.1 there). Note that $\mathbb{E} \left[P |\widehat{h}_n - \widehat{h}_{n-1}| \right] \leq \mathbb{E} P \widehat{h}_n^2 + \mathbb{E} P \widehat{h}_{n-1}^2$.

The next theorem gathers a slight refinement of classical deviation inequalities for the excess risk and the excess empirical risk. Complete derivations can be found in (Massart and Nédélec 2006; Massart 2007). A different style of presentation and derivation can be found in (Koltchinskii 2006). A roadmap of the proof is given in the appendix following the blend introduced in (Boucheron et al. 2005a).

Theorem 1 *Let \mathcal{H} denote the loss class associated with a learning task that satisfies the complexity Assumption (1) with function ψ and the noise conditions Assumption (2) with function ω . Let r_* denote the positive solution of*

$$\sqrt{n} r^2 = \psi(\omega(r)) .$$

Let h^* denote the relative loss function associated with the minimizer of the risk over all measurable functions. Then, with probability larger than $1 - 2\delta$, the excess risk $-P\hat{h}_n$ and the excess empirical risk $P_n\hat{h}_n$ satisfy:

$$\max\left(-P\hat{h}_n, P_n\hat{h}_n\right) \leq \kappa \left(Ph^* + r_*^2 + \frac{\omega^2(r_*)}{nr_*^2} \log \frac{1}{\delta} \right).$$

and

$$\max\left(\mathbb{E}\left[-P\hat{h}_n\right], \mathbb{E}\left[P_n\hat{h}_n\right]\right) \leq \kappa \left(Ph^* + r_*^2 \right).$$

In both inequalities, κ denotes a universal constant larger than 1 that does not depend on the learning task.

Remark 3 In the stochastic upper-bound, $\omega^2(r_*)/(nr_*^2)$ may be further upper-bounded by r_*^2 . This upper-bound may be used to simplify the statement of the Theorem, but it is important to keep the result in this form in order to derive meaningful results in Sections 4 and 5.

We may now combine the second inequality from Proposition 1 and the preceding Theorem to get meaningful upper-bounds for the variance of the empirical excess risk.

Theorem 2 *There exists a universal constant κ such that, using notations and assumptions of Theorem 1, the variance of the empirical excess risk is upper-bounded by*

$$\text{Var}\left[nP_n\hat{h}_n\right] \leq n\kappa \left(\omega^2(r_*) + \omega^2(\sqrt{Ph^*}) \right).$$

Proof Thanks to Proposition 1, we have:

$$\text{Var}\left[nP_n\hat{h}_n\right] \leq 2n\mathbb{E}\left[(P_{n-1} - P)\hat{h}_{n-1}\right] + 2n\mathbb{E}\left[P\hat{h}_n^2\right].$$

Thanks to the fact that $\omega^2(\sqrt{\cdot}) \in \mathcal{C}_1$, the last summand can be upper-bounded:

$$\begin{aligned} 2n\mathbb{E}\left[P\hat{h}_n^2\right] &\leq 4n\left(\mathbb{E}\left[P(\hat{h}_n - h^*)^2\right] + Ph^{*2}\right) \quad \text{as } (x+y)^2 \leq 2(x^2 + y^2) \\ &\leq 4n\left(\mathbb{E}\left[\omega^2\left(\sqrt{P(h^* - \hat{h}_n)}\right)\right] + \omega^2(\sqrt{Ph^*})\right) \\ &\leq 4n\left(2\omega^2\left(\sqrt{\mathbb{E}[P(h^* - \hat{h}_n)]}\right) + \omega^2(\sqrt{Ph^*})\right). \end{aligned}$$

Let us now relate $\mathbb{E}[(P_{n-1} - P)\hat{h}_{n-1}]$ with bounds on the excess empirical risk and the excess risk over samples of size n . Let s_* denote the positive solution of $\sqrt{n-1}r^2 = \psi(2\omega(r))$. It follows from Theorem 1 that:

$$\mathbb{E}\left[(P_{n-1} - P)\hat{h}_{n-1}\right] \leq 2\kappa \left(Ph^* + s_*^2 \right).$$

Now as $\sqrt{n-1}s_*^2 = \psi(\omega(s_*))$ and $\sqrt{nr_*^2} = \psi(\omega(r_*))$, we have $1 \leq s_*^2/r_*^2 \leq \frac{n}{n-1}$. So, plugging upper-bounds on expected and empirical excess risk from Theorem 1:

$$\begin{aligned} &\text{Var}\left[nP_n\hat{h}_n\right] \\ &\leq 4\kappa n \left(\frac{nr_*^2}{n-1} + Ph^* \right) + 8n\omega^2\left(\sqrt{\kappa(Ph^* + r_*^2)}\right) + 12n\omega^2(\sqrt{Ph^*}) \\ &\leq \kappa' n \left(r_*^2 + Ph^* + \omega^2(\sqrt{Ph^*}) + \omega^2(r_*) \right) \\ &\leq \kappa'' n \left(\omega^2(\sqrt{Ph^*}) + \omega^2(r_*) \right). \end{aligned}$$

□

Remark 4 The second bound in Proposition 1 provides an easy way to relate the variance of the empirical excess risk with r_* and the model bias Ph^* . With a little more work, the first bound would lead to the same kind of variance upper-bounds and it will prove very convenient when deriving higher moment estimates in the next section.

3.3 Application to bounded regression

Let us illustrate the variance bounds on a problem that provides a bridge between statistical learning and traditional statistics: bounded regression using least square estimation. This simple setting has been used to investigate model selection using (careful) data-dependent penalization by Arlot and Massart (2009). The latter paper takes advantage of the generalized Wilks phenomenon to build and calibrate nearly minimal penalties and to derive adaptive estimators.

Recall that the quadratic loss is defined by $\ell(y, y') = (y - y')^2$. Labels Y are assumed to be $[-1/2, 1/2]$ -valued. In the setting of the paper, the pseudo-distance d satisfying (1) and used in the statement of the complexity assumption and of the noise conditions assumption in Section 2 is defined in the following way. Let $f, f' \in \mathcal{F}$ then

$$d^2(f, f') = 2P\ell(f(X), f'(X)).$$

If $\bar{f} = f^*$ (we will see that this is a sensible assumption), the pseudo-distance can be related with the $L_2(P)$ norm of the (relative) loss functions:

$$\begin{aligned} Ph^2 &= P\left((f(X) - f^*(X))^2 ((f(X) - f^*(X)) - 2(Y - f^*(X)))^2\right) \\ &= P\left((f(X) - f^*(X))^2 \left((f(X) - f^*(X))^2 + 4(Y - f^*(X))^2\right)\right) \\ &\leq P(f(X) - f^*(X))^4 + 4P\left((f(X) - f^*(X))^2 (Y - f^*(X))^2\right) \\ &\leq 2P(f - f^*)^2 \\ &= d^2(f, f^*) \end{aligned}$$

where the second equation is obtained by reasoning conditionally on X , and the inequalities follow from the boundedness assumptions on Y , f and f' . As $R(f) - R(f^*) = P(f - f^*)^2$, we may choose $\omega(r) = \sqrt{2}r$ in order to comply with Assumption 2.

We consider a very simple estimator, the *regressogram*. The input space is partitioned into a finite (say D) number of pairwise disjoint subsets. Let Π denote the partition. The class \mathcal{F} is defined as the set of functions that are constant over the subsets defining Π . Any function f in \mathcal{F} is defined by a set of coefficients $(a_I)_{I \in \Pi}$:

$$f(x) = \sum_{I \in \Pi} a_I \mathbb{1}_I(x)$$

Given a sample of labelled points $(X_i, Y_i)_{i \leq n}$ where Y_i are $[-1/2, 1/2]$ -valued, the regressogram is the minimizer \hat{f} of the empirical quadratic risk over \mathcal{F} . It is defined by the set of coefficients $(\hat{a}_I)_{I \in \Pi}$

$$\hat{a}_I = \begin{cases} \frac{P_n(Y \mathbb{1}_I(X))}{P_n \mathbb{1}_I(X)} = \frac{P_n(Y \mathbb{1}_I(X))}{P_n(I)} & \text{if } P_n(I) > 0 \\ 0 & \text{otherwise.} \end{cases}$$

Note that the pseudo-distance is generated by a weighted pseudo-norm, if we let $f, f' \in \mathcal{F}$ be defined by vectors $\mathbf{a} = (a_I)_{I \in \Pi}$ and $\mathbf{b} = (b_I)_{I \in \Pi}$,

$$\|\mathbf{a}\|_2^2 = \sum_{I \in \Pi} P(I) a_I^2.$$

then $d^2(f, f') = 2\|\mathbf{a} - \mathbf{b}\|_2^2$.

Thanks to Pythagorean relationships, as \mathcal{F} is a closed subspace of $L_2(P)$, we may replace f^* by its projection onto \mathcal{F} . And thus, without loss of generality, we may assume that f^* belongs to \mathcal{F} . The function f^* is defined by the vector of coefficients \mathbf{a}^* defined by

$$a_I^* = \begin{cases} \frac{P(Y \mathbb{1}_I(X))}{P \mathbb{1}_I(X)} = \frac{P(Y \mathbb{1}_I(X))}{P(I)} & \text{if } P(I) > 0 \\ 0 & \text{otherwise.} \end{cases}$$

In order to simplify notations and computations, we further assume that f^* is the null function and that Y may be considered as additive homoschedastic centered noise. Let σ^2 denote the noise variance. Details concerning the general setting can be found in (Arlot and Massart 2009).

In that setting, we may choose

$$\psi(r) = 2r\sqrt{D}. \quad (7)$$

Indeed

$$(P_n - P)h = 2P_n(Yf(X)) - P_n f^2 + P f^2.$$

$$\begin{aligned} \mathbb{E} \left[\sup_{f: d^2(f, f^*) \leq r^2} |(P_n - P)h| \right] &\leq \mathbb{E} \left[\sup_{f: d^2(f, f^*) \leq r^2} |2P_n(Yf(X))| \right] \\ &\quad + \mathbb{E} \left[\sup_{f: d^2(f, f^*) \leq r^2} |(P - P_n)f^2| \right]. \end{aligned}$$

Recall that thanks to the assumption $f^* = 0$, $d^2(f, f^*) = 2\|\mathbf{a}\|^2$. Hence

$$\begin{aligned} \sup_{f: d^2(f, f^*) \leq r^2} |2P_n(Yf(X))| &= \frac{2}{n} \sup_{\mathbf{a}: 2\|\mathbf{a}\|^2 \leq r^2} \sum_{I \in \Pi} \sum_{i \leq n} (Y_i \mathbb{1}_I(X_i)) a_I \\ &\leq \frac{2}{n} \sup_{\mathbf{a}: \sqrt{2}\|\mathbf{a}\| \leq r} \|\mathbf{a}\| \sqrt{\sum_{I \in \Pi} \frac{(\sum_i Y_i \mathbb{1}_I(X_i))^2}{P(I)}} \\ &\quad \text{by Cauchy-Schwarz inequality.} \\ &\leq \frac{\sqrt{2}r}{n} \sqrt{\sum_{I \in \Pi} \frac{(\sum_i Y_i \mathbb{1}_I(X_i))^2}{P(I)}}. \end{aligned}$$

Taking expected values

$$\mathbb{E} \left[\sup_{f: d^2(f, f^*) \leq r^2} |2P_n(Yf(X))| \right] \leq \frac{\sqrt{2}r}{\sqrt{n}} \sqrt{D}.$$

On the other hand, using Cauchy-Schwarz inequality once again,

$$\mathbb{E} \left[\sup_{f: d^2(f, f^*) \leq r^2} |(P - P_n)f^2| \right] \leq \frac{r}{\sqrt{8}} \mathbb{E} \left[\left(\sum_{I \in \Pi} \frac{((P_n - P)(I))^2}{P(I)} \right)^{1/2} \right].$$

The right-hand side is the expected value of the square root of a Pearson statistics. Its expected value can be upper-bounded by $\sqrt{D/n}$. Inequality (7) follows from $\sqrt{2} + 1/\sqrt{8} \leq 2$.

Hence, r_* may be chosen as

$$r_* = \sqrt{\frac{8D}{n}}.$$

Hence, from Theorem 2, we obtain the following upper-bound on the variance of the unnormalized empirical excess risk:

$$8\kappa D.$$

In the simple setting considered here, it is possible to lower bound the expected value of the empirical excess risk. In order to alleviate computations, we further assume that $P(I) = 1/D$ for all $I \in \Pi$.

$$\begin{aligned} R_n(f^*) - R_n(\hat{f}) &= P_n \left(\sum_{I \in \Pi} (Y^2 - (Y - \hat{a}_I)^2) \mathbb{1}_I(X) \right) \\ &= \sum_{I \in \Pi} P_n(I) \hat{a}_I^2. \end{aligned}$$

Now, taking expectation over samples,

$$\begin{aligned} \mathbb{E} [R_n(f^*) - R_n(\hat{f})] &= \frac{\sigma^2}{n} \sum_{I \in \Pi} \mathbb{E}[\mathbb{1}_{P_n(I) > 0}] \\ &= \frac{\sigma^2}{n} \sum_{I \in \Pi} ((1 - (1 - P(I))^n)) \\ &\geq \frac{\sigma^2}{n} \sum_{I \in \Pi} (1 - \exp(-n/D)). \end{aligned}$$

If $n/D \geq \log(2)$, the last sum is larger than $D/2$. Note that $\mathbb{E} [R_n(f^*) - R_n(\hat{f})] \leq D\sigma^2/n$.

For a fixed value of σ^2 , the expected value and the variance of the unnormalized empirical excess risk are both proportional to D . As far as the first two moments are concerned, the bounded regression problem exemplifies the high-dimensional Wilks phenomenon.

4 Higher moments

The purpose of this section is to check that the excess empirical risk satisfies a Bernstein-like inequality. The classical Bernstein inequality for sums of independent random variables asserts the following. Let X_1, \dots, X_n be independent real-valued random variables. Assume that there exist positive numbers v and c such that $\sum_{i=1}^n \mathbb{E} [X_i^2] \leq v$ and

$$\sum_{i=1}^n \mathbb{E} [(X_i)_+^k] \leq \frac{k!}{2} v c^{k-2} \quad \text{for all integers } k \geq 3.$$

If $S = \sum_{i=1}^n (X_i - \mathbb{E}X_i)$, then for all $\lambda \in (0, 1/c)$ and $t > 0$,

$$\log \mathbb{E}e^{\lambda S} \leq \frac{v\lambda^2}{2(1-c\lambda)},$$

so that for all $t > 0$,

$$P\left\{S \geq \sqrt{2vt} + ct\right\} \leq e^{-t}.$$

Tail bounds, bounds on the logarithmic moment generating functions are known to be equivalent to moment bounds (see Ledoux and Talagrand 1991, for the ‘‘equivalence of moments principle’’). Let Y be a centered random variable. It can be checked that if for some positive v and c

$$P\left\{Y > \sqrt{2vt} + ct\right\} \wedge P\left\{-Y > \sqrt{2vt} + ct\right\} \leq e^{-t} \text{ for } t > 0,$$

then for every integer $q \geq 1$

$$\mathbb{E}[Y^{2q}] \leq q!(8v)^q + (2q!)(4c)^{2q}.$$

And conversely, if for some positive constants A, B , for every integer $q \geq 1$

$$\mathbb{E}[Y^{2q}] \leq q!A^q + (2q!)B^{2q}$$

then

$$P\left\{Y > \sqrt{8(A+B^2)t} + 2Bt\right\} \wedge P\left\{-Y > \sqrt{8(A+B^2)t} + 2Bt\right\} \leq e^{-t} \text{ for } t > 0.$$

Recall that if the X_i are recentered independent $\Gamma(v/c^2, 1/c)$ -distributed, then the preceding conditions are satisfied. This is true if the X_i are recentered χ_1^2 random variables, then $c = 2$.

Asserting that an M-estimation problem (for example a learning problem) exhibits a high dimensional Wilks phenomenon means that the q -norm of the centered empirical excess risk is upper-bounded by $\kappa(\sqrt{qv} + cq)$ where κ is a (another) universal constant, v is a sharp variance upper-bound and c depends slowly on sample size and model complexity/dimension, and is related to the ratio between the expected value and the variance.

Theorem 1 from the preceding section, provides deviation inequalities for empirical excess risk. Moreover, the statement of the theorem reveals that the upper-tail is not heavier than the upper-tail of a Gamma-shaped distribution. Those deviation inequalities can be completed by concentration inequalities that show that the excess empirical risk $P_n \hat{h}_n$ is sharply concentrated around its mean. Deviation inequalities provide information on the probability that a random variable deviates from its median or expectation by a large amount, they may be loose concerning the fluctuations of the random variable around its expectation. Concentration inequalities attempt to provide meaningful information on the probability of small as well as large fluctuations. The results of this section readily follow from the following moment inequalities (Boucheron et al. 2005b). Henceforth for $q > 0$, for any random variable X , $\|X\|_q^q = \mathbb{E}[|X|^q]$.

Theorem 3 *Let Z be a measurable function F of n independent random variables X_1, \dots, X_n , let X'_1, \dots, X'_n , be independent copies of X_1, \dots, X_n . Let $Z'_i = F(X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n)$. Let $V_+ = \sum_{i=1}^n \mathbb{E}'[(Z - Z'_i)_+^2]$ where \mathbb{E}' means conditional expectation conditioned on X_1, \dots, X_n . Then for any $q \geq 2$:*

$$\|(Z - \mathbb{E}[Z])_+\|_q \leq \sqrt{3q \|V_+\|_{q/2}} = \sqrt{3q} \left\| \sqrt{V_+} \right\|_q.$$

Let M be a random variable satisfying $(Z - Z'_i)_+ \leq M$ for all $i \leq n$, then for all $q \geq 2$

$$\|(Z - \mathbb{E}[Z])_-\|_q \leq \sqrt{5q} \left(\|\sqrt{V_+}\|_q \vee \|M\|_q \right).$$

In the sequel, just as in statement of Theorems 1 and 2, let \mathcal{H} denote the relative loss class associated with a learning task. From the preceding section, we know that if Z is the excess empirical risk, then the corresponding V_+ is upper-bounded by $2n(P_n \hat{h}_n^2 + P \hat{h}_n^2)$. Combining this observation with the theorem above immediately leads to the following proposition.

Proposition 2 *Let \mathcal{H} denote a separable class of functions. Let $Z = n \sup_{h \in \mathcal{H}} P_n h$ denote the unnormalized supremum of the empirical risk indexed by \mathcal{H} . Assume that the supremum is achieved at some $\hat{h}_n \in \mathcal{H}$, that is $Z = n P_n \hat{h}_n$. Then for $q \geq 2$*

$$\|(Z - \mathbb{E}[Z])_+\|_q \leq \sqrt{6nq} \left(\sqrt{\|P_n \hat{h}_n^2\|_{q/2}} + \sqrt{\|P \hat{h}_n^2\|_{q/2}} \right).$$

Proof Starting from

$$V_+ \leq 2n \left(P_n \hat{h}_n^2 + P \hat{h}_n^2 \right),$$

by Theorem 3, for $q \geq 2$:

$$\begin{aligned} & \|(Z - \mathbb{E}[Z])_+\|_q \\ & \leq \sqrt{3q} \left\| \sqrt{2n \left(P_n \hat{h}_n^2 + P \hat{h}_n^2 \right)} \right\|_q \\ & \leq \sqrt{6nq} \left(\sqrt{\|P_n \hat{h}_n^2\|_{q/2}} + \sqrt{\|P \hat{h}_n^2\|_{q/2}} \right). \end{aligned}$$

□

Note that the Talagrand Inequality for suprema of empirical processes (Talagrand 1996b) does provide us with a Bernstein-like inequality for the empirical excess risk $Z = n P_n \hat{h}_n$: for some universal constant κ ,

$$\|Z - \mathbb{E}[Z]\|_q \leq \kappa \left(\sqrt{q \mathbb{E} \left[\sup_{h \in \mathcal{H}} n P_n h^2 \right]} + q \right).$$

But, this inequality does not meet our requirements. In many situations, the variance upper-bounds it relies on, far exceed the upper-bounds described in the preceding section. Those bounds grow at least as fast as $n \sup_{h \in \mathcal{H}} P h^2$, that is linearly with sample size and cannot be easily related with model complexity that is with a sensible notion of degrees of freedom.

In the sequel, we assume that the learning task satisfies the complexity Assumption (1) with function ψ and the noise conditions Assumption (2) with function ω . Let r_* denote the positive solution of

$$\sqrt{nr}^2 = \psi(\omega(r)).$$

Recall that h^* denotes the relative loss function associated with the minimizer of the risk over all measurable functions, so $P h^*$ denotes the model bias. In the sequel, we will use the fact that with high probability, $P_n \hat{h}_n^2$ is not much larger than $P \hat{h}_n^2$. This is embodied in the following theorem whose proof can be found in the appendix.

Theorem 4 Let \mathcal{H} denote the loss class associated with a learning task that satisfies Assumption (1) about complexity with function ψ and Assumption (2) about noise conditions with function ω . Let r_* denote the positive solution of $\sqrt{nr}^2 = \psi(\omega(r))$. Let \hat{h}_n denote the relative loss function associated with the empirical risk minimizer. For $q \geq 2$

$$\begin{aligned} & \left\| P\hat{h}_n^2 \right\|_q \vee \left\| P_n\hat{h}_n^2 \right\|_q \\ & \leq \kappa \left(\omega^2 \left(\sqrt{Ph^*} \right) + \omega^2 \left(r_* \right) + \omega^2 \left(\frac{\omega(r_*)}{\sqrt{nr_*}} \right) q \right), \end{aligned}$$

where h^* denotes the relative loss function associated with the minimizer of the risk over all measurable functions.

The universal constant κ does not depend on the learning task.

Combining this theorem, the moment inequalities stated above and Proposition 2, leads to the main result:

Theorem 5 Let \mathcal{H} denote the loss class associated with a learning task that satisfies the complexity Assumption (1) with function ψ and the noise conditions Assumption (2) with function ω . Let r_* denote the positive solution of $\sqrt{nr}^2 = \psi(\omega(r))$. Let \hat{h}_n (resp. h^*) denote the relative loss function associated with the empirical risk minimizer (resp. the global minimizer of the risk). Let $Z = nP_n\hat{h}_n$ denote the unnormalized empirical excess risk. There exists a universal constant κ such that for $q \geq 2$,

$$\begin{aligned} & \|Z - \mathbb{E}[Z]\|_q \\ & \leq \kappa \left(\sqrt{n} \left(\omega \left(\sqrt{Ph^*} \right) + \omega(r_*) \right) q^{1/2} + \sqrt{n} \omega \left(\frac{\omega(r_*)}{\sqrt{nr_*}} \right) q \right). \end{aligned}$$

Remark 5 Comparing this inequality with the moment bounds that follow from the naive application of Talagrand's inequality shows that we are likely to obtain a significant improvement on the variance term which describes the sub-gaussian part of the tail. But we may feel less confident with respect to the second term on the right-hand-side. The latter describes the sub-exponential part of the tail. We will see in the next section that in some favorable situations $\sqrt{n} \omega \left(\frac{\omega(r_*)}{\sqrt{nr_*}} \right)$ can be upper-bounded by an expression that depends neither on sample size nor on model complexity but on noise conditions.

Proof The q -th norm of $P_n\hat{h}_n^2$ is upper-bounded due to Theorem 4. Now plugging the bounds from Theorem 4 into the last inequality leads to:

$$\begin{aligned} & \|(Z - \mathbb{E}[Z])_+\|_q \\ & \leq 2\sqrt{6n\kappa} \left(\omega \left(\sqrt{Ph^*} \right) + \omega(r_*) \right) q^{1/2} + 2\sqrt{3n\kappa} \omega \left(\frac{\omega(r_*)}{\sqrt{nr_*}} \right) q. \end{aligned}$$

The proof of the theorem is completed by observing that

$$\|(Z - \mathbb{E}[Z])\|_q \leq \|(Z - \mathbb{E}[Z])_+\|_q + \|(Z - \mathbb{E}[Z])_-\|_q$$

and that the random variable denoted by M in the statement of Theorem 3 is upper-bounded by 1. \square

We may turn again to the bounded regression problem which served as a benchmark for variance bounds at the end of Section 3. From Theorem 5, it follows that for a fixed noise level σ , there exists a universal constant κ such that the empirical excess risk satisfies:

$$\|Z - \mathbb{E}Z\|_q \leq \kappa(\sqrt{Dq} + q)$$

for all sample sizes n and all partition cardinalities D . On the other hand, under mild conditions on the relationship between the number of classes and the sample size, $\mathbb{E}Z/\sigma^2 \approx D$, so in this setting again, the empirical excess risk exemplifies the high dimensional Wilks phenomenon.

5 Application: classification with Vapnik-Chervonenkis classes

In this section, we plug the above-described moment estimates into the various bounds described by Massart and Nédélec (2006). Details concerning Vapnik-Chervonenkis classes can be found in (Assouad 1983; Vapnik. 1995; van der Vaart and Wellner 1996).

Proposition 3 *There exists a universal constant κ such that the following holds. Let \mathcal{F} denote a VC-class of sets with VC dimension V . Assume the following noise conditions: the Bayes classifier f^* belongs to \mathcal{F} , and conditionally on X , with probability $(1 + \beta)/2$, the label Y equals $f^*(X)$. The excess empirical risk satisfies*

$$\mathbb{E} \left[n(R_n(\bar{f}) - R_n(\hat{f})) \right] \leq \kappa \left[\left(\frac{V(1 + \log((n\beta^2/V) \vee 1))}{\beta} \right) \wedge \sqrt{nV} \right],$$

$$\text{Var} \left[n(R_n(\bar{f}) - R_n(\hat{f})) \right] \leq \kappa \left[\left(\frac{V(1 + \log((n\beta^2/V) \vee 1))}{\beta^2} \right) \wedge \sqrt{\frac{nV}{\beta^2}} \right]$$

and for $q \geq 2$

$$\begin{aligned} & \left\| \left(n(R_n(\bar{f}) - R_n(\hat{f})) - \mathbb{E}[n(R_n(\bar{f}) - R_n(\hat{f}))] \right)_+ \right\|_q \\ & \leq \sqrt{\kappa \left(\frac{V(1 + \log((n\beta^2/V) \vee 1))}{\beta^2} \right)} \wedge \sqrt{\frac{nV}{\beta^2}} q^{1/2} + \frac{\kappa}{\beta} q. \end{aligned}$$

Remark 6 These noise conditions have sometimes been called random classification noise in early Machine Learning papers, see Angluin and Laird (1987); Kearns et al. (1997); Bartlett et al. (2002).

Remark 7 For $n \geq V/\beta^2$, the upper-bounds stated in the proposition coincides with the one obtained when deriving tail bounds for a Gamma distribution with shape parameter $V(1 + \log((n\beta^2/V)))$ and scale parameter $1/\beta$. Up to the slowly varying function $\log((n\beta^2/V))$, this is in-line with what we meant using the expression "Wilks phenomenon".

The proof consists in putting together the generic moment bounds described in this paper and the detailed computations carried out in (Massart and Nédélec 2006).

Proof In the random classification noise setting, for any classifier f

$$R(f) - R(f^*) = \beta \mathbb{E} [\mathbb{1}_{f \neq f^*}] = \beta \mathbb{E} \left[(f(X) - f^*(X))^2 \right].$$

Hence the function ω may be chosen as $\omega(r) = 1 \vee r/\sqrt{\beta}$. On the other hand, computations carried out in (Massart and Nédélec 2006) show that ψ may be chosen as

$$\psi(r) = Kr\sqrt{V(1 + \log(r^{-1} \vee 1))}$$

where K is some universal constant that depends neither on the VC-class under investigation nor on the sampling probability distribution. The quantity denoted by r_* is the solution of

$$\sqrt{nr}^2 = K \frac{r}{\sqrt{\beta}} \sqrt{V \left((1 + \log \left(\frac{\sqrt{\beta}}{r} \vee 1 \right)) \right)}.$$

A few lines of algebra (see Massart and Nédélec (2006)) lead to

$$r_*^2 \leq K^2 \left[\left(\frac{V(1 + \log((n\beta^2/V) \vee 1))}{n\beta} \right) \wedge \sqrt{\frac{V}{n}} \right]$$

so that Theorem 1 and Proposition 2 lead to the first two inequalities.

Note that $\frac{\omega(r_*)}{\sqrt{nr_*}} = \frac{1}{\sqrt{n\beta}}$. Hence the bounds on the q -th moment of the excess empirical risk translate into the last inequality. \square

Acknowledgements The authors would like to thank Sylvain Arlot for helpful comments and for pointing a mistake in a previous draft.

References

- H. Akaike. A new look at the statistical model identification. *IEEE Trans. on Automatic Control*, 19(6): 716–723, 1974.
- K. Alexander. Rates of growth and sample moduli for weighted empirical processes indexed by sets. *Probab. Theory and Related Fields*, 75:379–423, 1987.
- D. Angluin and P. Laird. Learning from noisy examples. *Machine Learning*, 2(4):343–370, 1987.
- S. Arlot and P. Massart. Data-driven calibration of penalties for least-squares regression. *Journal of Machine Learning Research*, 10:245–279, 2009.
- P. Assouad. Densité et dimension. *Ann. Inst. Fourier (Grenoble)*, 33(3):233–282, 1983.
- P. Bartlett and S. Mendelson. Empirical minimization. *Probab. Theory Related Fields*, 135(3):311–334, 2006.
- P. Bartlett, S. Boucheron, and G. Lugosi. Model selection and error estimation. *Machine Learning*, 48: 85–113, 2002.
- P. Bickel and K. Doksum. *Mathematical statistics*. Holden-Day Inc., San Francisco, Calif., 1976.
- S. Boucheron, O. Bousquet, and G. Lugosi. Theory of classification: some recent advances. *ESAIM Probability & Statistics*, 9:329–375, 2005a.
- S. Boucheron, O. Bousquet, G. Lugosi, and P. Massart. Moment inequalities for functions of independent random variables. *Annals of Probability*, 33(2):514–560, 2005b.
- O. Bousquet. Concentration inequalities for sub-additive functions using the entropy method. In *Stochastic inequalities and applications*, volume 56 of *Progr. Probab.*, pages 213–247. Birkhäuser, Basel, 2003.
- O. Bousquet. A Bennett concentration inequality and its application to suprema of empirical processes. *C. R. Math. Acad. Sci. Paris*, 334(6):495–500, 2002.
- V. de la Pena and E. Giné. *Decoupling*. Springer-Verlag, 1999.
- L. Devroye and T. Wagner. Distribution-free inequalities for the deleted and holdout error estimates. *IEEE Trans. Inform. Theory*, 25:202–207, 1977.
- B. Efron and C. Stein. The jackknife estimate of variance. *Annals of Statistics*, 9(3):586–596, 1981.

- J. Fan. Local linear regression smoothers and their minimax efficiency. *Annals of Statistics*, 21:196–216, 1993.
- J. Fan, C. Zhang, and J. Zhang. Generalized likelihood ratio statistics and Wilks phenomenon. *Annals of Statistics*, 29(1):153–193, 2001.
- G. Gayraud and C. Pouet. Minimax Testing Composite Null Hypotheses in the Discrete Regression Scheme. *Mathematical Methods of Statistics*, 10(4):375–394, 2001.
- E. Giné and V. Koltchinskii. Concentration inequalities and asymptotic results for ratio type empirical processes. *Annals of Probability*, 34(3):1143–1216, 2006.
- E. Giné, V. Koltchinskii, and J. Wellner. *Stochastic Inequalities and Applications*, chapter Ratio limit theorems for empirical processes, pages 249–278. Birkhäuser, 2003.
- P. Huber. The behavior of the maximum likelihood estimates under non-standard conditions. In *Proc. Fifth Berkeley Symposium on Probability and Mathematical Statistics*, pages 221–233. Univ. California Press, 1967.
- M. Kearns and D. Ron. Algorithmic stability and sanity-check bounds for leave-one-out cross-validation. *Neural Computation*, 11(6):1427–1453, 1999.
- M. Kearns, Y. Mansour, A. Ng, and D. Ron. An experimental and theoretical comparison of model selection methods. *Machine Learning*, 27:7–50, 1997.
- V. Koltchinskii. Localized rademacher complexities and oracle inequalities in risk minimization. *Annals of Statistics*, 34:2593–2656, 2006.
- M. Ledoux. On Talagrand’s deviation inequalities for product measures. *ESAIM Probability & Statistics*, 1: 63–87 (electronic), 1995/97.
- M. Ledoux. *The concentration of measure phenomenon*. AMS, 2001.
- M. Ledoux and M. Talagrand. *Probability in Banach spaces*. Springer, 1991.
- C. Mallows. Some comments on C_p . *Technometrics*, 15(4):661–675, 1973.
- E. Mammen and A. Tsybakov. Smooth discrimination analysis. *Annals of Statistics*, 27(6):1808–1829, 1999.
- P. Massart. About the constants in Talagrand’s concentration inequality. *Annals of Probability*, 28:863–885, 2000a.
- P. Massart. Some applications of concentration inequalities to statistics. *Ann. Fac. Sc. Toulouse*, IX(2): 245–303, 2000b.
- P. Massart. *Concentration inequalities and model selection*. *Ecole d’Été de Probabilité de Saint-Flour XXXIV*, volume 1896 of *Lecture Notes in Mathematics*. Springer-Verlag, 2007.
- P. Massart and E. Nédélec. Risk bounds for classification. *Annals of Statistics*, 34(5):2326, 2006.
- D. Pollard. *Convergence of stochastic processes*. Springer-Verlag, 1984.
- S. Portnoy. Asymptotic behavior of likelihood methods for exponential families when the number of parameters tends to infinity. *Annals of Statistics*, 16:356–366, 1988.
- M. Quenouille. Approximate test of correlation in time series. *Journal of the Royal Statistical Society, Ser. B*, 11:68–84, 1949.
- A. Rakhlin, S. Mukherjee, and T. Poggio. Stability Results in Learning Theory. *Analysis and Applications (Singapore)*, 3:397–417, 2005.
- E. Rio. Inégalités de concentration pour les processus empiriques de classes de parties. *Probab. Theory Related Fields*, 119:163–175, 2001.
- B. Schoelkopf and A. Smola. *Learning with kernels*. MIT Press, 2002.
- G. R. Shorack and J. A. Wellner. *Empirical processes with applications to statistics*. John Wiley & Sons Inc., 1986.
- M. Talagrand. A new look at independence. *Annals of Probability*, 24:1–34, 1996a.
- M. Talagrand. New concentration inequalities in product spaces. *Inventiones Mathematicae*, 126:505–563, 1996b.
- A. B. Tsybakov. Optimal aggregation of classifiers in statistical learning. *Annals of Statistics*, 32:135–166, 2004.
- J. Tukey. Bias and confidence in not quite large samples. *Annals of Mathematical Statistics*, 29:614, 1958.
- S. van de Geer. *Applications of empirical process theory*. Cambridge University Press, 2000.
- A. van der Vaart. *Asymptotic statistics*. Cambridge University Press, 1998.
- A. van der Vaart and J. Wellner. *Weak convergence and empirical processes*. Springer-Verlag, 1996.
- V. Vapnik. *The nature of statistical learning theory*. Springer-Verlag, 1995.
- S. Wilks. The large-sample distribution of the likelihood ratio for testing composite hypotheses. *Annals of Mathematical Statistics*, 9:60–62, 1938.

A Tools from empirical process theory

This appendix provides a sketch of proof of Theorems 1 and 4. The proof of these two results depends on the following concentration inequality for suprema of empirical processes. It is due to Talagrand (Talagrand 1996b, Theorem 1.4). The sharp constants in the last inequality are due to Bousquet (2002) building on previous work by several authors (Ledoux 1995/97; Massart 2000a; Rio 2001).

Theorem 6 *There exists a number κ with the following property. Consider n independent random variables U_1, \dots, U_n valued in some measurable space Ω . Consider a countable class \mathcal{H} of measurable functions on Ω . Consider the random variable $Z = \sup_{h \in \mathcal{H}} \sum_{i=1}^n h(U_i)$. Consider $u = \sup_{h \in \mathcal{H}, x \in \Omega} |h(x)|$ and $v = \mathbb{E} \sup_{h \in \mathcal{H}} \sum_{i=1}^n h(X_i)^2$. Then for each $t > 0$*

$$\Pr \{|Z - \mathbb{E}Z| > t\} \leq K \exp \left(-\frac{1}{K} \frac{t}{u} \log \left(1 + \frac{t}{uv} \right) \right).$$

If $P_h = 0$ for all $h \in \mathcal{H}$ and $u = 1$, then letting $\sigma^2 = \sup_{h \in \mathcal{H}} Ph^2$ and $v = 2\mathbb{E}[Z] + n\sigma^2$, for all $\delta > 0$,

$$\mathbb{P} \left\{ Z \geq \mathbb{E}[Z] + \sqrt{2v \log \frac{1}{\delta}} + \frac{1}{3} \log \frac{1}{\delta} \right\} \leq \delta.$$

The development of tight bounds for excess risk and empirical excess risk relies on a good understanding of the stochastic behavior of the increments of the empirical process indexed by \mathcal{F} . Indeed, the sum of the excess risk and the empirical excess risk (which are both positive) coincides with the increment of the centered empirical process defined by the loss functions between \bar{f} and \hat{f} :

$$R(\hat{f}) - R(\bar{f}) + R_n(\bar{f}) - R_n(\hat{f}) = (P - P_n) \left(\ell(\hat{f}(X), Y) - \ell(\bar{f}(X), Y) \right) = (P_n - P)\hat{h}.$$

If we are given a (random) function ϕ_n such that for all classifiers f in \mathcal{F} ,

$$|(P - P_n) (\ell(f(X), Y) - \ell(\bar{f}(X), Y))| \leq \phi_n(R(f)),$$

a sharp upper-bound on $R(\hat{f}) - R(\bar{f})$ is obtained by solving the inequality $R(f) - R(\bar{f}) \leq \phi_n(R(f))$. The analysis of the modulus of oscillation of empirical processes indexed by general classes of sets goes back to the work of Alexander (1987) (see also van de Geer 2000, for applications to M -estimation), and the impact of recently obtained concentration inequalities on this topic is thoroughly investigated by Giné et al. (2003) and Giné and Koltchinskii (2006). However, the analysis of the excess risk does not require an exhaustive understanding of the modulus of continuity but a tight control of this modulus over a special range of values. Hence, rather than resorting to omnibus results like those presented in the aforementioned references, we rely on ad hoc analysis which has already proved convenient and successful (Massart 2000b; Massart and Nédélec 2006; Massart 2007).

The following lemma is borrowed from (Massart and Nédélec 2006).

Lemma 1 *Let \mathcal{H} be some countable index set. Let L denote a non-negative function on \mathcal{H} , and assume there exists some $\bar{h} \in \mathcal{H}$ such that $L(\bar{h}) = \inf_{h \in \mathcal{H}} L(h)$.*

Let Z denote a stochastic process indexed by \mathcal{H} . Assume that there exists a function ψ from the class \mathcal{C}_1 of positive sub-linear functions such that for all r satisfying $r^2 \geq \psi(r)$:

$$\mathbb{E} \left[\sup_{h: h \in \mathcal{H}, L(h) \leq r^2} |Z(h) - Z(\bar{h})| \right] \leq \psi(r).$$

Then, for any r satisfying $r^2 \geq \psi(r)$:

$$\mathbb{E} \left[\sup_{h \in \mathcal{H}} \frac{r^2}{r^2 + L(h)} |Z(h) - Z(\bar{h})| \right] \leq 4\psi(r).$$

The next theorem is the key technical ingredient in the derivation of tail bounds for excess risk and excess empirical risk (Theorem 1) as well as for the $L_2(P)$ and $L_2(P_n)$ distances of relative loss \hat{h}_n to 0 (Theorem 4).

Theorem 7 Let \mathcal{H} denote a countable class of measurable functions over \mathcal{X} . Let L denote a non-negative function on \mathcal{H} , which achieves its minimum at $0 \in \mathcal{H}$. Assume furthermore that $\sup_{x \in \mathcal{X}} |h(x)| \leq 1$.

Let σ be some non-negative mapping on \mathcal{H} such that for every $h \in \mathcal{H}$ $Ph^2 \leq \sigma^2(h)$. Assume furthermore that there exists some function ρ from \mathcal{C}_1 such that for any $h \in \mathcal{H}$: $\sigma(h) \leq \rho(\sqrt{L(h)})$. Assume that there exists some function ψ from \mathcal{C}_1 such that for r satisfying $\psi(r) \leq \sqrt{nr^2}$

$$\sqrt{n}\mathbb{E} \left[\sup_{h \in \mathcal{H}, \sigma(h) \leq r} |(P_n - P)h| \right] \leq \psi(r).$$

Let K denote a number larger than 1. Let δ and ϵ be positive numbers. Assume that the positive solution $r(\delta)$ of equation

$$r^2 = K(1 + \epsilon) \left(4 \frac{\psi(\rho(r))}{\sqrt{n}} + \rho(r) \sqrt{\frac{\log \frac{1}{\delta}}{n} + \frac{1}{\epsilon} \frac{\log \frac{1}{\delta}}{n}} \right),$$

satisfies $\psi(r) \leq \sqrt{nr^2}$. Then, with probability larger than $1 - 2\delta$, for all $h \in \mathcal{H}$

$$|(P - P_n)h| \leq \frac{L(h) + r^2(\delta)}{K}.$$

Proof We prove that with probability larger than $1 - \delta$, for all $h \in \mathcal{H}$

$$(P_n - P)h \leq \frac{L(h) + r^2(\delta)}{K}.$$

Let r satisfy $\psi(r) \leq \sqrt{nr^2}$, and let the random variable V_r be defined as

$$V_r = \sup_{h \in \mathcal{H}} r^2 \frac{(P_n - P)h}{L(h) + r^2}.$$

Then V_r is the supremum of a centered empirical process indexed by \mathcal{H} . Moreover

$$\mathbb{E} \left[\sup_{h, L(h) \leq r^2} |(P_n - P)h| \right] \leq \mathbb{E} \left[\sup_{\sigma(h) \leq \rho(r)} |(P_n - P)h| \right] \leq \frac{\psi(\rho(r))}{\sqrt{n}}.$$

As $\psi \circ \rho \in \mathcal{C}_1$, by Lemma 1, we have $\mathbb{E}[V_r] \leq 4 \frac{\psi(\rho(r))}{\sqrt{n}}$. Now, note that if $L(h) \leq r^2$,

$$\text{Var} \left[\frac{r^2 h}{L(h) + r^2} \right] \leq Ph^2 \leq \rho^2(r).$$

Moreover, for each $h \in \mathcal{F}$,

$$\text{Var} \left[\frac{r^2 h}{L(h) + r^2} \right] \leq \left(\frac{r^2 \rho(\sqrt{L(h)})}{L(h) + r^2} \right)^2 \leq \left(\frac{r^2 \rho(\sqrt{L(h)})}{2\sqrt{L(h)}r} \right)^2 \leq \left(\frac{r \rho(\sqrt{L(h)})}{2\sqrt{L(h)}} \right)^2.$$

Thus if $L(h) \geq r^2$, as ρ is assumed to belong to the class \mathcal{C}_1 (of positive sub-linear functions),

$$\text{Var} \left[\frac{r^2 h}{L(h) + r^2} \right] \leq \frac{\rho^2(r)}{4}.$$

On the other hand, for all $h \in \mathcal{H}$

$$\sup_{x \in \mathcal{X}} \left| \frac{r^2 h(x)}{L(h) + r^2} \right| \leq 1.$$

We may now invoke Talagrand's inequality. With probability larger than $1 - \delta$ for all $h \in \mathcal{H}$

$$\begin{aligned} V_r &\leq \mathbb{E}[V_r] + \sqrt{\frac{2}{n} (2\mathbb{E}[V_r] + \rho^2(r)) \log \frac{1}{\delta}} + \frac{1}{3n} \log \frac{1}{\delta} \\ &\leq (1 + \epsilon)\mathbb{E}[V_r] + \rho(r) \sqrt{\frac{\log \frac{1}{\delta}}{n}} + \left(\frac{1}{3} + \frac{1}{\epsilon} \right) \frac{\log \frac{1}{\delta}}{n} \\ &\leq 4(1 + \epsilon) \frac{\psi(\rho(r))}{\sqrt{n}} + \rho(r) \sqrt{\frac{\log \frac{1}{\delta}}{n}} + \left(\frac{1}{3} + \frac{1}{\epsilon} \right) \frac{\log \frac{1}{\delta}}{n}. \end{aligned}$$

Let $r(\delta)$ be defined as in the statement of the Theorem. With probability larger than $1 - \delta$, for all $h \in \mathcal{H}$:

$$(P_n - P)h \leq \frac{1}{K} (L(h) + r^2(\delta)) .$$

The same line of reasoning can be applied to the process $((P - P_n)h)_{h \in \mathcal{H}}$. This leads to the theorem. \square

The following corollary reveals that if the neighborhood of 0 in \mathcal{H} is not too rich, the empirical $L_2(P_n)$ structure of \mathcal{H} around 0 faithfully approximates the $L_2(P)$ structure of \mathcal{H} .

Corollary 1 *Let \mathcal{H} denote a countable class of measurable functions over \mathcal{X} . And assume that $\sup_{x \in \mathcal{X}} |h(x)| \leq 1$. Assume that there exists some function $\phi \in \mathcal{C}_1$ such that for all r such that $\sqrt{nr^2} \geq \phi(r)$:*

$$\sqrt{n} \mathbb{E} \left[\sup_{h \in \mathcal{H}, Ph^2 \leq r^2} |(P_n - P)h| \right] \leq \phi(r) .$$

Let K be larger than 1, let ϵ be a positive number. Assume that be the positive solution $r(\delta)$ of equation

$$r^2 = K(1 + \epsilon) \left(\frac{32\phi(r)}{\sqrt{n}} + r \sqrt{\frac{\log \frac{1}{\delta}}{n}} + \frac{1}{\epsilon} \frac{\log \frac{1}{\delta}}{n} \right)$$

satisfies $\sqrt{nr^2}(\delta) \geq \phi(r(\delta))$. Then, with probability larger than $1 - 2\delta$, for all $h \in \mathcal{H}$

$$-\frac{1}{K}(Ph^2 + r^2(\delta)) \leq (P_n - P)h^2 \leq \frac{1}{K}(Ph^2 + r^2(\delta)) .$$

Proof The proof of Corollary 1 consists in applying Theorem 7 to the class of functions $\{h^2 : h \in \mathcal{H}\}$. Note first that

$$\sup_{h \in \mathcal{H}, x \in \mathcal{X}} |h^2(x) - Ph^2| \leq 1 .$$

Second, choosing $L(h) = \sigma^2(h) = Ph^2$, L is minimized by 0, and letting $\rho = \text{Id}$,

$$P(h^2)^2 \leq \sigma^2(h) \leq \rho^2(\sqrt{L(h)}) .$$

The only point that needs to be checked is that

$$\sqrt{n} \mathbb{E} \left[\sup_{h: \sigma(h) \leq r} |(P_n - P)h^2| \right] \leq 8\phi(r) .$$

Let now $\epsilon_1, \dots, \epsilon_n$ be a sequence of independent symmetric, $\{-1, 1\}$ -valued random variables, the so-called Rademacher variables. By a standard symmetrization argument:

$$\mathbb{E} \left[\sup_{h: \sigma(h) \leq r} |(P_n - P)h^2| \right] \leq \frac{2}{n} \mathbb{E} \left[\mathbb{E} \left[\sup_{h: \sigma(h) \leq r} \left| \sum_{i=1}^n \epsilon_i h^2(X_i) \right| \mid X_1, \dots, X_n \right] \right]$$

where the internal expectation is taken over the Rademacher variables.

Now as $x^1 \rightarrow x^2$ is 2-Lipschitz over $[-1, 1]$, the contraction principle (Ledoux and Talagrand 1991; de la Pena and Giné 1999) implies that

$$\frac{1}{n} \mathbb{E} \left[\sup_{h: \sigma(h) \leq r} \left| \sum_{i=1}^n \epsilon_i h^2(X_i) \right| \mid X_1, \dots, X_n \right] \leq \frac{2}{n} \mathbb{E} \left[\sup_{h: \sigma(h) \leq r} \left| \sum_{i=1}^n \epsilon_i h(X_i) \right| \mid X_1, \dots, X_n \right] .$$

Another standard symmetrization argument entails

$$\mathbb{E} \left[\sup_{h: \sigma(h) \leq r} |(P_n - P)h^2| \right] \leq 8 \mathbb{E} \left[\sup_{h: \sigma(h) \leq r} |(P - P_n)h| \right] \leq 8\phi(r) .$$

We may now invoke Theorem 7 on the class $\{h^2 : h \in \mathcal{H}\}$, with $L(h) = \sigma^2(h) = Ph^2$, $\rho = \text{Id}$ and $\psi = 8\phi$. \square

The following elementary proposition allows us to translate stochastic bounds in the style of (Koltchinskii 2006) into bounds stated in the style of (Massart and Nédélec 2006). It provides the ultimate ingredient in the proof of Theorems 1 and 4.

Proposition 4 *Let ψ and ρ denote non-constant, functions from \mathcal{C}_1 . Let r_* denote the unique positive solution of equation $\sqrt{nr^2} = \psi(\rho(r))$. For some a, b, c in \mathbb{R}_+ , with $a \geq 1$, let s denote the unique solution of equation $r^2 = \frac{a}{\sqrt{n}}\psi(\rho(r)) + \frac{b}{\sqrt{n}}\rho(r) + c$. Then $s^2 \leq 2a^2r_*^2 + 2b^2\frac{\rho^2(r_*)}{nr_*^2} + 2c$.*

Proof (Theorem 1) The proof consists of invoking Theorem 7 for the relative loss class \mathcal{H} defined as on page 6. From the definitions of \mathcal{H} and $L(h) = P(h^* - h)$, L is minimized by 0. Moreover, the function σ and ρ can be chosen as $\sigma(h) = d(f, f^*)$ and $\rho = 2\omega$. Indeed,

$$\begin{aligned} Ph^2 &\leq 2(P(h^* - h)^2 + P(h^*)^2) \\ &\leq 2(d^2(f, f^*) + d^2(\bar{f}, f^*)) \\ &\leq 2\left(\omega^2(\sqrt{P(h^* - h)}) + \omega^2(\sqrt{Ph^*})\right) \\ &\leq 4\omega^2(\sqrt{P(h^* - h)}) = 4\omega^2(\sqrt{L(h)}). \end{aligned}$$

We may now invoke Theorem 7 to check that, with probability larger than $1 - 2\delta$, for all $h \in \mathcal{H}$:

$$(P - P_n)h \leq \frac{1}{K} (L(h) + r^2(\delta)).$$

Going back to classifiers and risks, this translates into

$$\frac{K-1}{K}(-Ph) + P_n h \leq \frac{1}{K} (P(h^* - h) + r^2(\delta)) = \frac{1}{K} (Ph^* + r^2(\delta)).$$

Specializing the preceding bound to \hat{h}_n , we can take advantage of the fact that both $-P\hat{h}_n$ and $P_n\hat{h}_n$ are both non-negative. With probability larger $1 - 2\delta$, the excess risk is upper-bounded by

$$-P\hat{h}_n \leq \frac{1}{K-1} (Ph^* + r^2(\delta)),$$

while the empirical excess risk is upper-bounded by

$$P_n\hat{h}_n \leq \frac{1}{K} (Ph^* + r^2(\delta)).$$

The proof can be completed by invoking Proposition 4 and integration by parts. \square

Proof (Theorem 4.) From the proof of Theorem 1, it follows that with probability larger than $1 - 2\delta$

$$L(\hat{h}_n) \leq \frac{KL(0)}{K-1} + \frac{r^2(\delta)}{K-1}.$$

Then, from the relation

$$P\hat{h}_n^2 \leq 4\omega^2 \left(\sqrt{L(\hat{h}_n)} \right).$$

and thanks to the fact that $\omega^2(\sqrt{\cdot}) \in \mathcal{C}_1$, with probability $1 - 2\delta$

$$\frac{1}{4}P\hat{h}_n^2 \leq \frac{K}{K-1} \left(\omega^2 \left(\sqrt{Ph^*} \right) + \omega^2 \left(\sqrt{\frac{r^2(\delta)}{K}} \right) \right).$$

Now, invoking the relation between $r^2(\delta)$ and r_*^2 (Proposition 4), this last inequality translates into

$$\frac{1}{4}P\hat{h}_n^2 \leq \frac{K}{K-1} \left(\omega^2 \left(\sqrt{Ph^*} \right) + 2Ka^2\omega^2(r_*) + 2K\omega^2 \left(\frac{\omega(r_*)}{\sqrt{nr_*}} \right) \log \frac{1}{\delta} \right). \quad (8)$$

The inequality concerning $P\hat{h}_n^2$ follows by integration by parts.

The inequality concerning $P_n \widehat{h}^2$ follows by combining the inequality concerning $P \widehat{h}_n^2$ and the stochastic relation between the $L_2(P)$ and the $L_2(P_n)$ structure of the loss class \mathcal{H} (Corollary 1). Note that if we keep using K and δ , and call $s(\delta)$ the solution of

$$r^2 = K \left(8a \frac{\phi(r)}{\sqrt{n}} + r \sqrt{\frac{\log \frac{1}{\delta}}{n}} + \frac{c \log \frac{1}{\delta}}{n} \right),$$

then $s(\delta) \leq r(\delta)$ since $2\omega(r) \geq r$ on $[0, 1]$. This entails that with probability $1 - 4\delta$, we have both (8) and

$$P_n \widehat{h}_n^2 \leq \frac{K+1}{K-1} P \widehat{h}_n^2 + \frac{r^2(\delta)}{K}.$$

The Theorem follows by observing that $r(\delta) \geq r_*$ and $\omega(r) \geq r$ on $[0, 1]$, and using Proposition 4. \square