



# DFG-SNF Research Group FOR916

Statistical Regularization and Qualitative Constraints

Sara van de Geer

Johannes Lederer

## The Bernstein-Orlicz norm and deviation inequalities

Preprint FOR916 11-6

Updated Version (20-11-2011)

Preprint-Series of the Research Group FOR916

# The Bernstein-Orlicz norm and deviation inequalities

Sara van de Geer · Johannes Lederer

Received: date / Accepted: date

**Abstract** We introduce two new concepts designed for the study of empirical processes. First, we introduce a new Orlicz norm which we call the Bernstein-Orlicz norm. This new norm interpolates sub-Gaussian and sub-exponential tail behavior. In particular, we show how this norm can be used to simplify the derivation of deviation inequalities for suprema of collections of random variables. Secondly, we introduce chaining and generic chaining along a tree. These simplify the well-known concepts of chaining and generic chaining. The supremum of the empirical process is then studied as a special case. We show that chaining along a tree can be done using entropy with bracketing. Finally, we establish a deviation inequality for the empirical process for the unbounded case.

**Keywords** Bernstein's inequality · Chaining along a tree · Deviation inequality · Empirical process · Orlicz norm

**Mathematics Subject Classification (2000)** 60E15 · 60F10

## 1 Introduction

We introduce a new Orlicz norm which we name the Bernstein-Orlicz norm. It interpolates sub-Gaussian and sub-exponential tail behavior. With this new norm, we apply the usual techniques based on Orlicz norms. In particular, we derive deviation inequalities for suprema in a fairly simple and straightforward way. The Bernstein-Orlicz norm captures Bernstein's probability inequalities, and its use puts further

---

Research supported by SNF 20PA21-120050.

S. van de Geer  
Seminar for Statistics, ETH Zürich  
Rämistrasse 101, 8092 Zürich, Switzerland  
Tel.: 00-41-44-6322252  
Fax: 00-41-44-6321228  
E-mail: geer@stat.math.ethz.ch

J. Lederer  
Seminar for Statistics, ETH Zürich  
Rämistrasse 101, 8092 Zürich, Switzerland

derivations in a unifying framework, shared for example by techniques for the sub-Gaussian case, such as those for empirical processes based on symmetrization and Hoeffding's inequality.

We furthermore introduce chaining and generic chaining along a tree, which is we believe conceptually simpler than the usual chaining and generic chaining. We invoke it for the presentation of maximal inequalities for general random variables with finite Bernstein-Orlicz norm. The supremum of the empirical process is then studied as a special case, and we show that chaining along a tree can be done using entropy with bracketing. We establish a deviation inequality for the empirical process indexed by a class of functions  $\mathcal{G}$ , in terms of the new Bernstein-Orlicz norm. The class  $\mathcal{G}$  is assumed to satisfy a uniform Bernstein condition, but need not be uniformly bounded in supremum norm.

The paper is organized as follows. In Section 2, we introduce the Bernstein-Orlicz norm and discuss the relation with Bernstein's inequality. We then present some bounds for maxima of finitely many random variables (Section 3) or suprema over a countable set of random variables (Section 4). Section 4 also contains the concept of (generic) chaining along a tree. The proofs of the results in Sections 2, 3 and 4 are elementary and given immediately following their statement. Section 5 contains the application to the empirical process. The proofs here are more technical, and given separately in Sections 6 and 7.

## 2 The Bernstein-Orlicz norm

Consider a random variable  $Z \in \mathbb{R}$  with distribution  $\mathbb{P}$ . We first recall the general Orlicz norm (see e.g. Krasnosel'skii and Rutickii [1961]).

**Definition 1** Let  $\Psi : [0, \infty) \rightarrow [0, \infty)$  be an increasing and convex function with  $\Psi(0) = 0$ . The  $\Psi$ -Orlicz norm of  $Z$  is

$$\|Z\|_{\Psi} := \inf \left\{ c > 0 : \mathbb{E} \Psi \left( \frac{|Z|}{c} \right) \leq 1 \right\}.$$

A special case is the  $L_m(\mathbb{P})$ -norm ( $m \geq 1$ ) which corresponds to  $\Psi(z) = z^m$ . Other important special cases are  $\Psi(z) = \exp[z^2] - 1$  for sub-Gaussian random variables and  $\Psi(z) = \exp(z) - 1$  for sub-exponential random variables. We propose functions  $\Psi$  that combine sub-Gaussian intermediate tails and sub-exponential far tails.

For each  $L > 0$  we define

$$\Psi_L(z) := \exp \left[ \frac{\sqrt{1 + 2Lz} - 1}{L} \right]^2 - 1, \quad z \geq 0. \quad (1)$$

It is easy to see that  $\Psi_L$  is increasing and convex, and that  $\Psi_L(0) = 0$ .

**Definition 2** Let  $L > 0$  be given. The ( $L$ -)Bernstein-Orlicz norm is the  $\Psi$ -Orlicz norm with  $\Psi = \Psi_L$  given in (1).

Indeed, the Bernstein-Orlicz norm combines sub-Gaussian and sub-exponential behavior:

$$\Psi_L(z) \approx \begin{cases} \exp[z^2] - 1 & \text{for } Lz \text{ small} \\ \exp[2z/L] - 1 & \text{for } Lz \text{ large} \end{cases}.$$

Note that the constant  $L$  governs the range of the sub-Gaussian behavior. It is a dimensionless constant, i.e., it does not depend on the scale of measurement.

The inverse of  $\Psi_L$  is

$$\Psi_L^{-1}(t) = \sqrt{\log(1+t)} + \frac{L}{2} \log(1+t), \quad t \geq 0.$$

With this and with Chebyshev's inequality, one now directly derives a probability inequality for  $Z$ .

**Lemma 1** *Let  $\tau := \|Z\|_{\Psi_L}$ . We have for all  $t > 0$ ,*

$$\mathbb{P}\left(|Z| > \tau \left[\sqrt{t} + \frac{Lt}{2}\right]\right) \leq 2 \exp[-t].$$

**Proof of Lemma 1.** By Chebyshev's inequality, for all  $c > \|Z\|_{\Psi_L}$ ,

$$\begin{aligned} \mathbb{P}\left(|Z|/c \geq \sqrt{t} + \frac{Lt}{2}\right) &= \mathbb{P}\left(|Z|/c \geq \Psi_L^{-1}(e^t - 1)\right) \\ &= \mathbb{P}\left(\Psi_L(|Z|/c) \geq e^t - 1\right) \leq \left(\mathbb{E}\Psi_L(|Z|/c) + 1\right) e^{-t}. \end{aligned}$$

Thus,

$$\begin{aligned} \mathbb{P}\left(|Z|/\tau > \sqrt{t} + \frac{Lt}{2}\right) &= \lim_{c \downarrow \tau} \mathbb{P}\left(|Z|/c > \sqrt{t} + \frac{Lt}{2}\right) \\ &\leq \lim_{c \downarrow \tau} \left(\mathbb{E}\Psi_L(|Z|/c) + 1\right) e^{-t} \leq 2e^{-t}. \end{aligned}$$

□

The next lemma says that a converse result holds as well, that is, from the probability inequality of Lemma 1 one can derive a bound for the Bernstein-Orlicz norm, with constants  $L$  and  $\tau$  multiplied by  $\sqrt{3}^1$ .

**Lemma 2** *Suppose that for some constants  $\tau$  and  $L$ , and for all  $t > 0$ ,*

$$\mathbb{P}\left(|Z| \geq \tau \left[\sqrt{t} + \frac{Lt}{2}\right]\right) \leq 2 \exp[-t].$$

*Then  $\|Z\|_{\Psi_{\sqrt{3}L}} \leq \sqrt{3}\tau$ .*

**Proof of Lemma 2.** We have

$$\begin{aligned} \mathbb{E}\Psi_{\sqrt{3}L}\left(|Z|/(\sqrt{3}\tau)\right) &= \int_0^\infty \mathbb{P}\left(|Z| \geq \sqrt{3}\tau \Psi_{\sqrt{3}L}^{-1}(t)\right) dt \\ &= \int_0^\infty \mathbb{P}\left(|Z| \geq \sqrt{3}\tau \left[\sqrt{\log(1+t)} + \frac{\sqrt{3}L}{2} \log(1+t)\right]\right) dt \\ &= \int_0^\infty \mathbb{P}\left(|Z| \geq \tau \left[\sqrt{\log(1+t)^3} + \frac{L}{2} \log(1+t)^3\right]\right) dt \leq 2 \int_0^\infty \frac{1}{(1+t)^3} dt = 1. \end{aligned}$$

□

We recall Bernstein's inequality, see Bennet [1962].

<sup>1</sup> The constant can possibly be improved.

**Theorem 1** Let  $X_1, \dots, X_n$  be independent random variables with values in  $\mathbb{R}$  and with mean zero. Suppose that for some constants  $\sigma$  and  $K$ , one has

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}|X_i|^m \leq \frac{m!}{2} K^{m-2} \sigma^2, \quad m = 1, 2, \dots$$

Then for all  $t > 0$ ,

$$\mathbb{P}\left(\frac{1}{\sqrt{n}} \left| \sum_{i=1}^n X_i \right| \geq \sigma\sqrt{2t} + \frac{Kt}{\sqrt{n}}\right) \leq 2 \exp[-t].$$

The following corollary shows that  $\|\cdot\|_{\Psi_L}$  indeed captures the nature of Bernstein's inequality.

**Corollary 1** Let  $X_1, \dots, X_n$  be independent random variables satisfying the conditions of Theorem 1. Then by this theorem and Lemma 2, for  $L := \sqrt{6}K/(\sqrt{n}\sigma)$ , we have

$$\left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \right\|_{\Psi_L} \leq \sqrt{6}\sigma.$$

### 3 The Bernstein-Orlicz norm for the maximum of finitely many variables

Using Orlicz norms, the argument for obtaining a bound for the expectation of maxima is standard. We refer to van der Vaart and Wellner [1996] for a general approach. We consider the special case of the Bernstein-Orlicz norm.

**Lemma 3** Let  $\tau$  and  $L$  be constants, and let  $Z_1, \dots, Z_p$  be random variables satisfying

$$\max_{1 \leq j \leq p} \|Z_j\|_{\Psi_L} \leq \tau.$$

Then

$$\mathbb{E} \max_{1 \leq j \leq p} |Z_j| \leq \tau \left[ \sqrt{\log(1+p)} + \frac{L}{2} \log(1+p) \right].$$

**Proof of Lemma 3 .** Let  $c > \tau$ . Then by Jensen's inequality

$$\begin{aligned} \mathbb{E} \max_{1 \leq j \leq p} |Z_j| &\leq c\Psi_L^{-1} \left( \mathbb{E} \Psi_L \left( \max_{1 \leq j \leq p} |Z_j|/c \right) \right) = c\Psi_L^{-1} \left( \mathbb{E} \max_{1 \leq j \leq p} \Psi_L \left( |Z_j|/c \right) \right) \\ &\leq c\Psi_L^{-1} \left( \sum_{j=1}^p \mathbb{E} \Psi_L \left( |Z_j|/c \right) \right) \leq c\Psi_L^{-1} \left( p \max_{1 \leq j \leq p} \mathbb{E} \Psi_L \left( |Z_j|/c \right) \right). \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbb{E} \max_{1 \leq j \leq p} |Z_j| &\leq \lim_{c \downarrow \tau} c\Psi_L^{-1} \left( p \max_{1 \leq j \leq p} \mathbb{E} \Psi_L \left( |Z_j|/c \right) \right) \leq \tau\Psi_L^{-1}(p) \\ &= \tau \left[ \sqrt{\log(1+p)} + \frac{L}{2} \log(1+p) \right]. \end{aligned}$$

□

As a special case, one may consider the random variables

$$Z_j := \frac{1}{\sqrt{n}} \sum_{i=1}^n g_j(X_i), \quad j = 1, \dots, p,$$

where  $X_1, \dots, X_n$  are independent random variables with values in some space  $\mathcal{X}$ , and where  $g_1, \dots, g_p$  are real-valued functions on  $\mathcal{X}$ . If the  $g_j(X_i)$  are centered for all  $i$  and  $j$ , and if one assumes the Bernstein condition

$$\frac{1}{n} \sum_{i=1}^n \mathbf{E}|g_j(X_i)|^m \leq \frac{m!}{2} K^{m-2} \sigma^2, \quad m = 2, 3, \dots, \quad j = 1, \dots, p,$$

then one can apply Lemma 3, with  $\tau := \sqrt{6}\sigma$  and  $L = \sqrt{6}K/(\sqrt{n}\sigma)$ , giving the inequality

$$\mathbf{E} \max_{1 \leq j \leq p} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n g_j(X_i) \right| \leq \sigma \sqrt{6 \log(1+p)} + \frac{3K}{\sqrt{n}} \log(1+p). \quad (2)$$

This follows from Corollary 1. The constants can however be improved when using direct arguments (see e.g. Lemma 14.12 Bühlmann and van de Geer [2011]).

We now present a deviation inequality in probability for the maximum of finitely many variables.

**Lemma 4** *Let  $Z_1, \dots, Z_p$  be random variables satisfying for some  $L$  and  $\tau$*

$$\max_{1 \leq j \leq p} \|Z_j\|_{\Psi_L} \leq \tau.$$

*Then for all  $t > 0$*

$$\mathbf{P} \left( \max_{1 \leq j \leq p} |Z_j| \geq \tau \left[ \sqrt{\log(1+p)} + \frac{L}{2} \log(1+p) + \sqrt{t} + \frac{Lt}{2} \right] \right) \leq 2 \exp[-t].$$

**Proof of Lemma 4.** We first use that for any  $a > 0$  and  $t > 0$ , one has  $\sqrt{a} + \sqrt{t} > \sqrt{a+t}$ , so that

$$\begin{aligned} & \mathbf{P} \left( \max_{1 \leq j \leq p} |Z_j| \geq \tau \left[ \sqrt{\log(1+p)} + \frac{L}{2} \log(1+p) + \sqrt{t} + \frac{Lt}{2} \right] \right) \\ & \leq \mathbf{P} \left( \max_{1 \leq j \leq p} |Z_j| > \tau \left[ \sqrt{t + \log(1+p)} + \frac{L}{2} (t + \log(1+p)) \right] \right). \end{aligned}$$

Next, we apply the union bound and Lemma 1:

$$\begin{aligned} & \mathbf{P} \left( \max_{1 \leq j \leq p} |Z_j| > \tau \left[ \sqrt{t + \log(1+p)} + \frac{L}{2} (t + \log(1+p)) \right] \right) \\ & \leq \sum_{j=1}^p \mathbf{P} \left( |Z_j| > \tau \left[ \sqrt{t + \log(1+p)} + \frac{L}{2} (t + \log(1+p)) \right] \right) \\ & \leq 2p \exp \left[ -(t + \log(1+p)) \right] = \frac{2p}{1+p} \exp[-t] \leq 2 \exp[-t]. \end{aligned}$$

□

Using Lemma 2, this is easily converted into the following deviation inequality for the Bernstein-Orlicz norm. We use the notation

$$x_+ := x \mathbf{1}\{x > 0\}.$$

**Lemma 5** Let  $Z_1, \dots, Z_p$  be random variables satisfying for some  $L$  and  $\tau$

$$\max_{1 \leq j \leq p} \|Z_j\|_{\psi_L} \leq \tau.$$

Then

$$\left\| \left( \max_{1 \leq j \leq p} |Z_j| - \tau \left[ \sqrt{\log(1+p)} + \frac{L}{2} \log(1+p) \right] \right)_+ \right\|_{\psi_{\sqrt{3}L}} \leq \sqrt{3}\tau.$$

**Proof of Lemma 5.** Let

$$Z := \left( \max_{1 \leq j \leq p} |Z_j| - \tau \left[ \sqrt{\log(1+p)} + \frac{L}{2} \log(1+p) \right] \right)_+.$$

By Lemma 4, we have for all  $t > 0$

$$\begin{aligned} & \mathbb{P} \left( Z \geq \tau \left[ \sqrt{t} + \frac{Lt}{2} \right] \right) \\ &= \mathbb{P} \left( \max_{1 \leq j \leq p} |Z_j| \geq \tau \left[ \sqrt{\log(1+p)} + \frac{L}{2} \log(1+p) + \sqrt{t} + \frac{Lt}{2} \right] \right) \leq 2 \exp[-t]. \end{aligned}$$

Application of Lemma 2 finishes the proof. □.

#### 4 Chaining along a tree

A common technique for bounding suprema of stochastic processes is chaining as developed by Kolmogorov, leading to versions of Dudley's entropy bound (Dudley [1967]). See e.g. van der Vaart and Wellner [1996] or van de Geer [2000] and the references therein. We however propose another method which we call chaining along a tree. This method is conceptually simpler than the usual chaining and, as far as we know, does not introduce unnecessary restrictions. An example will be detailed in Section 5 for the case of entropy with bracketing. The generic chaining technique of Talagrand [2005] is a refinement which we shall also consider in Definition 6 and Theorem 3.

Let  $S \in \mathbb{N}_0$  be fixed.

**Definition 3** A finite tree<sup>2</sup>  $\mathcal{T}$  is a collection  $\{G_s\}_{s=0}^S$  of disjoint subsets of  $\{1, \dots, N\}$  such that  $\cup_{s=0}^S G_s = \{1, \dots, N\}$ , together with a function

$$\text{parent} : \{1, \dots, N\} \rightarrow \{1, \dots, N\},$$

such that  $\text{parent}(j) \in G_{s-1}$  for  $j \in G_s$ ,  $s \in 1, \dots, S$ . We call an element of  $\{1, \dots, N\}$  a node, and  $G_s$  a generation,  $s = 0, \dots, S$ . A branch of the tree with end node  $j_S \in G_S$  is the sequence  $\{j_0, \dots, j_S\}$  with  $j_{s-1} = \text{parent}(j_s)$ ,  $s = 1, \dots, S$ .

**Definition 4** Let a collection of real-valued random variables  $\mathcal{W} := \{W_j\}_{j=1}^N$  be given. A finite labeled tree  $(\mathcal{T}, \mathcal{W})$  is a finite tree with on each node  $j$  a label  $W_j$ .

Let  $\Theta$  be some countable set and let  $Z_\theta \in \mathbb{R}$  be a random variable defined for each  $\theta \in \Theta$ . We consider supremum of the process  $\{Z_\theta : \theta \in \Theta\}$ .

<sup>2</sup> Actually,  $\mathcal{T}$  is rather a forest consisting of  $|G_0|$  trees

**Definition 5** Let  $\delta > 0$  and  $\tau > 0$  be constants and let  $\mathcal{L} := \{L_s\}_{s=0}^S$  be a sequence of positive numbers. A  $(\delta, \tau, \mathcal{L})$  finite tree chain for  $\{Z_\theta\}$  is a finite labeled tree  $(\mathcal{T}, \mathcal{W})$  such that for all  $s = 0, \dots, S$ ,

$$\|W_j\|_{\Psi_{L_s}} \leq \tau 2^{-s}, \quad \forall j \in G_s,$$

and such that one can apply chaining of  $\{Z_\theta\}$  along the tree  $(\mathcal{T}, \mathcal{W})$ , with approximation error  $\delta$ . That is, for each  $\theta \in \Theta$  there is an end node  $j_S \in G_S$  such that the branch  $\{j_0, \dots, j_S\}$  satisfies

$$|Z_\theta| \leq \sum_{s=0}^S |W_{j_s}| + \delta.$$

In the above definition, the approximation error  $\delta$  will generally depend on the depth  $S$  of the tree. We assume that at a fine enough level, the approximation error is small. The usual chaining technique does not assume a tree structure, but indeed often needs only a finite number of steps. A tree structure follows if the members at the finest level are taken as end nodes. With a finite number of steps, the sum given in (3) is finite. This avoids requiring convergence of an infinite sum.

We have presented the definition of a finite tree chain for the Bernstein-Orlicz norm  $\|\cdot\|_{\Psi_L}$ . However, the concept is not particularly tied up with this norm, e.g., for sub-Gaussian cases one may choose to replace the Bernstein-Orlicz norm by the  $L_2(\mathbb{P})$  norm (corresponding to case where the constants in  $\mathcal{L}$  all vanish).

Let us now turn to the results.

**Theorem 2** Let  $(\mathcal{T}, \mathcal{W})$  be an  $(\delta, \tau, \mathcal{L})$  finite tree chain for  $\{Z_\theta\}$ . Define

$$\gamma := \tau \sum_{s=0}^S 2^{-s} \left[ \sqrt{\log(1 + |G_s|)} + \frac{L_s}{2} \log(1 + |G_s|) \right]. \quad (3)$$

It holds that

$$\mathbf{E} \left( \sup_{\theta \in \Theta} |Z_\theta| \right) \leq \gamma + \delta. \quad (4)$$

*Remark 1* One may minimize the right hand side of (4) over all finite trees.

**Proof of Theorem 2.** We have

$$\mathbf{E} \sup_{\theta \in \Theta} |Z_\theta| \leq \sum_{s=0}^S \mathbf{E} \max_{j \in G_s} |W_j| + \delta.$$

Application of Lemma 3 gives that for each  $s \in \{0, \dots, S\}$

$$\mathbf{E} \max_{j \in G_s} |W_j| \leq \tau 2^{-s} \left[ \sqrt{\log(1 + |G_s|)} + \frac{L_s}{2} \log(1 + |G_s|) \right].$$

□

With generic chaining, the condition on the Bernstein-Orlicz norm of the labels is dropped in the definition of the tree. This Bernstein-Orlicz norm then turns up in the constants (5) and (6) which appear in the generic chaining bound of Theorem 3

**Definition 6** Let  $\delta > 0$  be a constant. A  $\delta$  finite generic tree chain for  $\{Z_\theta\}$  is a finite labeled tree  $(\mathcal{T}, \mathcal{W})$  such that one can apply generic chaining of  $\{Z_\theta\}$  along the tree  $(\mathcal{T}, \mathcal{W})$  with approximation error  $\delta$ . That is, for each  $\theta \in \Theta$  there is an end node  $j_S \in G_S$  such that the branch  $\{j_0, \dots, j_S\}$  satisfies

$$|Z_\theta| \leq \sum_{s=0}^S |W_{j_s}| + \delta.$$

Let  $(\mathcal{T}, \mathcal{W})$  a finite labeled tree. For each end node  $k \in G_S$ , we let

$$\{j_0(k), \dots, j_S(k)\}$$

be the corresponding branch (so that  $j_S(k) = k$ ), and we write

$$W_s(k) := W_{j_s(k)}, \quad k \in G_S, \quad s = 0, 1, \dots, S.$$

Fix a sequence of positive constants  $\mathcal{L} := \{L_s\}_{s=0}^S$ . We write for  $k \in G_S$ ,

$$\gamma_{1,*}(k) := \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} \sqrt{\log(1 + |G_s|)}, \quad (5)$$

$$\gamma_{2,*}(k) := \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} L_s \log(1 + |G_s|), \quad (6)$$

$$\gamma_*(k) := \gamma_{1,*}(k) + \frac{\gamma_{2,*}(k)}{2}.$$

Moreover, we let

$$\gamma_{1,*} := \max_{k \in G_S} \gamma_{1,*}(k), \quad \gamma_{2,*} := \max_{k \in G_S} \gamma_{2,*}(k), \quad \gamma_* := \max_{k \in G_S} \gamma_*(k),$$

and

$$\tau_* := \max_{k \in G_S} \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} \sqrt{1 + s},$$

and

$$L_* \tau_* := \max_{k \in G_S} \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} (1 + s) L_s.$$

**Theorem 3** Let  $(\mathcal{T}, \mathcal{W})$  be a  $\delta$  finite generic tree chain for  $\{Z_\theta\}$ . Then

$$\mathbb{P} \left( \sup_{\theta \in \Theta} |Z_\theta| \geq \gamma_* + \delta + \tau_* \left[ 1 + \frac{L_*}{2} \right] + \tau_* \left[ \sqrt{t} + \frac{L_* t}{2} \right] \right) \leq 2 \exp[-t].$$

*Remark 2* The result of Theorem 3 may again be optimized over all finite generic trees.

**Proof of Theorem 3.** Define for  $s = 0, \dots, S$ ,

$$\alpha_s := \left[ \sqrt{\log(1 + |G_s|)} + \frac{L_s}{2} \log(1 + |G_s|) \right] + \left[ \sqrt{(1 + s)(1 + t)} + \frac{(1 + s)(1 + t)L_s}{2} \right].$$

Using Lemma 4, we see that

$$\mathbb{P} \left( \max_{j \in G_s} \frac{|W_j|}{\|W_j\|_{\Psi_{L_s}}} \geq \alpha_s \right) \leq 2 \exp[-(1 + t)(1 + s)], \quad s = 0, \dots, S. \quad (7)$$

We have

$$\begin{aligned}
& \mathbb{P}\left(\max_k \sum_{s=0}^S |W_s(k)| \geq \max_k \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} \alpha_s\right) \\
& \leq \mathbb{P}\left(\exists k : \sum_{s=0}^S |W_s(k)| \geq \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} \alpha_s\right) \\
& \leq \sum_{s=0}^S \mathbb{P}\left(\exists k : |W_s(k)| \geq \|W_s(k)\|_{\Psi_{L_s}} \alpha_s\right) \\
& = \sum_{s=0}^S \mathbb{P}\left(\max_k \frac{|W_s(k)|}{\|W_s(k)\|_{\Psi_{L_s}}} \geq \alpha_s\right) \leq \sum_{s=0}^S \mathbb{P}\left(\max_{j \in G_s} \frac{|W_j|}{\|W_j\|_{\Psi_{L_s}}} \geq \alpha_s\right).
\end{aligned}$$

Now insert (7) to find

$$\begin{aligned}
\mathbb{P}\left(\max_k \sum_{s=0}^S |W_s(k)| \geq \max_k \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} \alpha_s\right) & \leq 2 \sum_{s=0}^S \exp[-(1+t)(1+s)] \\
& \leq \frac{2e^{-(1+t)}}{1-e^{-(1+t)}} \leq \frac{2e^{-1}}{1-e^{-1}} \exp[-t] \leq 2 \exp[-t].
\end{aligned}$$

We have by definition

$$\begin{aligned}
\max_k \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} \left[ \sqrt{\log(1+|G_s|)} + \frac{L_s}{2} \log(1+|G_s|) \right] & = \gamma_*, \\
\max_k \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} \sqrt{(1+s)} & = \tau_*,
\end{aligned}$$

and

$$\max_k \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} (1+s)L_s = \tau_* L_*.$$

Therefore,

$$\begin{aligned}
\max_k \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} \alpha_s & \leq \gamma_* + \tau_* \sqrt{1+t} + \frac{\tau_*(1+t)L}{2} \\
& \leq \gamma_* + \tau_* + \frac{\tau_* L_*}{2} + \tau_* \left[ \sqrt{t} + \frac{L_* t}{2} \right].
\end{aligned}$$

□

Note that the constants  $L_*$  and  $\tau_*$  possibly depend on the complexity of  $\Theta$  through the quantities  $\{\|W_s(k)\|_{\Psi_{L_s}} : k \in G_s, s = 0, \dots, S\}$ . Moreover, the choice of the constants  $\mathcal{L} = \{L_s\}_{s=0}^S$  may also depend on the complexity of  $\Theta$ . In the application to the empirical process (see Section 5), the latter will be indeed the case. We will nevertheless derive there a deviation inequality where we put the dependency on the complexity of  $\Theta$  in the shift.

As a simple corollary of Theorem 3, one obtains a deviation inequality in the Bernstein-Orlicz norm. We state this for completeness. In Section 5 we will not apply Corollary 2 directly, because as such, it does not allow us to put all dependency on the complexity of  $\Theta$  in the shift.

**Corollary 2** *Let the conditions of Theorem 3 be met. Then the combination of this theorem with Lemma 2 gives*

$$\left\| \left( \sup_{\theta \in \Theta} |Z_\theta| - (\gamma_* + \delta + \tau_* [1 + L_*/2]) \right)_+ \right\|_{\psi_{\sqrt{3}L_*}} \leq \sqrt{3}\tau_*.$$

By Jensen's inequality, we then get

$$\mathbf{E} \sup_{\theta \in \Theta} |Z_\theta| \leq \gamma_* + \delta + \tau_* \left[ 1 + \frac{L_*}{2} \right] + \sqrt{3}\tau_* \left[ \sqrt{\log 2} + \frac{\sqrt{3}L_*}{2} \log 2 \right].$$

*Example 1* In Talagrand [2005], the sizes  $|G_s|$  of generation  $s$  is fixed to be

$$|G_s| = 2^{2^s}, \quad s = 0, \dots, S.$$

In that case,

$$\log(1 + |G_s|) \leq (2^{2^s} + 1) \log 2 \leq 2^{2^s+1} \leq 2^{2^{s+1}}.$$

Hence

$$\gamma_* \leq 2\gamma_0,$$

where

$$\gamma_0 := \max_{k \in G_S} \gamma_0(k),$$

and for  $k \in G_S$ ,

$$\gamma_0(k) := \gamma_{1,0}(k) + \frac{\gamma_{2,0}(k)}{2},$$

and

$$\gamma_{1,0}(k) := \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} 2^s, \quad \gamma_{2,0}(k) := \sum_{s=0}^S \|W_s(k)\|_{\Psi_{L_s}} L_s 2^{2s}.$$

Furthermore, since  $1 + s \leq 2^{2^s}$  for all  $s \geq 0$ ,

$$\tau_* \leq \gamma_{1,0} := \max_{k \in G_S} \gamma_{1,0}(k),$$

and

$$\tau_* L_* \leq \gamma_{2,0} := \max_{k \in G_S} \gamma_{2,0}(k).$$

Hence,

$$\gamma_* + \tau_* \left[ 1 + \frac{L_*}{2} \right] \leq 3 \left[ \gamma_{1,0} + \frac{\gamma_{2,0}}{2} \right],$$

and

$$\sqrt{3}\tau_* \left[ \sqrt{\log 2} + \sqrt{\frac{3L_*}{2}} \log 2 \right] \leq \sqrt{3 \log 2} \gamma_{1,0} + \frac{3 \log 2}{2} \gamma_{2,0}.$$

It follows from Corollary 2 that

$$\mathbf{E} \sup_{\theta \in \Theta} |Z_\theta| \leq (3 + \sqrt{3 \log 2}) \gamma_{1,0} + \frac{3 + 3 \log 2}{2} \gamma_{2,0}.$$

Thus, we arrive at a special case of Theorem 1.2.7 in Talagrand [2005]. The latter book does not treat deviation inequalities.

When using a  $(\delta, \tau, \mathcal{L})$  finite tree chain, one takes  $\|W_s(k)\|_{\Psi_{L_s}} \leq \tau 2^{-s}$  for all  $s$  and  $k \in G_s$ . In that case, the constants  $\tau_*$  and  $L_*$  in the bounds given in Corollary 2 only depend on the scale parameter  $\tau$  and on the constants  $\mathcal{L} = \{L_s\}_{s=0}^S$ . This is detailed in the next theorem.

**Theorem 4** *Let the conditions of Theorem 2 be met, and define*

$$\gamma := \tau \sum_{s=0}^S 2^{-s} \left[ \sqrt{\log(1 + |G_s|)} + \frac{L_s}{2} \log(1 + |G_s|) \right],$$

and

$$L := \sum_{s=0}^S \frac{2^{-s} L_s (1+s)}{4}.$$

Then for all  $t > 0$

$$\mathbf{P} \left( \sup_{\theta \in \Theta} |Z_\theta| \geq \gamma + \delta + 4\tau \left[ 1 + \frac{L}{2} \right] + 4\tau \left[ \sqrt{t} + \frac{Lt}{2} \right] \right) \leq 2 \exp[-t].$$

**Proof of Theorem 4.** This follows from Theorem 3, where one takes

$$\|W_s(k)\|_{\Psi_{L_s}} \leq \tau 2^{-s}.$$

We have

$$\tau_*/\tau \leq \sum_{s=0}^S 2^{-s} \sqrt{(1+s)} = 2 \sum_{s=1}^{S+1} 2^{-s} \sqrt{s} \leq 2 \int_0^\infty 2^{-x} \sqrt{x} dx = \frac{\sqrt{\pi}}{(\log 2)^{3/2}} \leq 4.$$

Moreover,

$$L_* \tau_* \leq \sum_{s=0}^S 2^{-s} L_s (1+s) = 4L.$$

□

## 5 Application to empirical processes

Let  $\mathcal{X}$  be some measurable space, and consider independent  $\mathcal{X}$ -valued random variables  $X_1, \dots, X_n$ . Let  $\mathcal{G}$  be a collection of real-valued functions on  $\mathcal{X}$ .

Write

$$P_n g := \frac{1}{n} \sum_{i=1}^n g(X_i), \quad P g := \frac{1}{n} \sum_{i=1}^n \mathbf{E} g(X_i),$$

and

$$\|g\|^2 := \frac{1}{n} \sum_{i=1}^n \mathbf{E} g^2(X_i).$$

We assume the normalization

$$\sup_{g \in \mathcal{G}} \|g\| \leq 1.$$

We study the supremum of the empirical process  $\{\nu_n(g) : g \in \mathcal{G}\}$ , where  $\nu_n(g) := \sqrt{n}(P_n - P)g$ .

We recall the deviation inequality of Massart [2000], which refines the constants in Talagrand [1996].

**Theorem 5** (Massart [2000]) Suppose that for a constant  $K$

$$\sup_{g \in \mathcal{G}} \sup_{x \in \mathcal{X}} |g(x)| \leq K. \quad (8)$$

Then for all  $\epsilon > 0$  and all  $t > 0$ , it holds that

$$\mathbb{P} \left( \sup_{g \in \mathcal{G}} |\nu_n(g)| \geq (1 + \epsilon) \mathbb{E} \sup_{g \in \mathcal{G}} |\nu_n(g)| + \sqrt{2\kappa t} + \kappa(\epsilon) K t / \sqrt{n} \right) \leq \exp[-t], \quad (9)$$

where  $\kappa$  and  $\kappa(\epsilon)$  can be taken equal to  $\kappa = 4$  and  $\kappa(\epsilon) = 2.5 + 32/\epsilon$ .

For the i.i.d. case, Bousquet [2002] obtained constants remarkably close those to for the case where  $\mathcal{G}$  is a singleton. In fact, Massart [2000] and Bousquet [2002] and others have derived *concentration* inequalities which in addition to upper bounds show similar lower bounds for the supremum of the empirical process. This is complemented in Lederer and van de Geer [2011] to moment concentration inequalities assuming only moment conditions on the envelope  $\Gamma(\cdot) := \sup_{g \in \mathcal{G}} |g(\cdot)|$ , instead of the boundedness assumption (8).

In this paper, we provide a deviation inequality of the same spirit as in the above Theorem 5, where we replace condition (8) by a weaker Bernstein condition (see (11)), which essentially requires that the  $g(X_i)$  have sub-exponential tails, and where we also present a deviation result in Bernstein-Orlicz norm. These deviation results in probability and in Bernstein-Orlicz norm are given in Theorem 8. We have not tried to optimize the constants. Moreover, we replace the expectation  $\mathbb{E} \sup_{g \in \mathcal{G}} |\nu_n(g)|$  in (9) by the upper bound we obtain from chaining arguments<sup>3</sup>. Deviation inequalities for the sub-exponential case can be found in literature (see e.g. Viens and Vizcarra [2007]), but these do not cover the more refined interpolation of sub-Gaussian and sub-exponential tail behavior. The above cited work also contains lower bounds for suprema, thus completing the results to concentration inequalities.

Now our first aim is to show that entropy with bracketing conditions allow one to construct a finite tree chain. We recall here the definition of a bracketing set and entropy with bracketing (see Blum [1955], or see van der Vaart and Wellner [1996], van de Geer [2000] and their references).

**Definition 7** Let  $s > 0$  be arbitrary. A  $2^{-s}$ -bracketing set for  $\{\mathcal{G}, \|\cdot\|\}$  is a finite collection of functions  $\{\{\tilde{g}_j^L, \tilde{g}_j^U\}\}_{j=1}^{\tilde{N}_s}$  satisfying  $\|\tilde{g}_j^U - \tilde{g}_j^L\| \leq 2^{-s}$  for all  $j$ , and such that for each  $g \in \mathcal{G}$  there is a  $j \in \{1, \dots, \tilde{N}_s\}$  such that  $\tilde{g}_j^L \leq g \leq \tilde{g}_j^U$ . If no such finite collection exists, we write  $\tilde{N}_s = \infty$ .

We also introduce a generalized bracketing set, in the spirit of van de Geer [2000].

**Definition 8** Let  $K > 0$  be a fixed constant. A generalized bracketing set for  $\mathcal{G}$  is a finite collection of functions  $\{\{\tilde{g}_j^L, \tilde{g}_j^U\}\}_{j=1}^{\tilde{N}_0}$  satisfying for all  $j$

$$P|\tilde{g}_j^U - \tilde{g}_j^L|^m \leq \frac{m!}{2} (2K)^{m-2}, \quad m = 2, 3, \dots,$$

and such that for each  $g \in \mathcal{G}$  there is a  $j \in \{1, \dots, \tilde{N}_0\}$  such that  $\tilde{g}_j^L \leq g \leq \tilde{g}_j^U$ . Write  $\tilde{N}_0 = \infty$  if no such finite collection exists.

<sup>3</sup> This upper bound can be shown to be (up to constants) tight in certain examples. The upper bound following from generic chaining is modulo constants tight for the general sub-Gaussian case.

A special case is where the envelope function  $\Gamma := \sup_{g \in \mathcal{G}} |g|$  satisfies the Bernstein condition

$$P\Gamma^m \leq \frac{m!}{2}(2K)^{m-2}, \quad m = 2, 3, \dots$$

Then one can take  $[-\Gamma, \Gamma]$  as generalized bracketing set, consisting of only one element.

In what follows, we let for each  $s \in \mathbb{N}$ ,  $\tilde{N}_s$  be the cardinality of a minimal  $2^{-s}$ -bracketing set for  $\mathcal{G}$ . The  $2^{-s}$ -entropy with bracketing of  $\mathcal{G}$  is

$$\tilde{H}_s := \log(1 + \tilde{N}_s), \quad s \in \mathbb{N}.$$

Moreover,  $\tilde{N}_0$  is the cardinality of a minimal generalized bracketing set, and we let

$$\tilde{H}_0 := \log(1 + \tilde{N}_0).$$

Finally, we write

$$N_s := \prod_{k=0}^s \tilde{N}_k, \quad H_s := \log(1 + N_s), \quad s \in \mathbb{N}_0. \quad (10)$$

The following theorem uses arguments of Ossiander [1987], and is comparable to Theorem 2.7.11 in Talagrand [2005] (who adapts the technique of Ossiander [1987]). However, we do not use generic chaining here. On the other hand, our results lead to the more involved deviation inequalities as given in Theorem 8.

**Theorem 6** *Suppose that for some constant  $K \geq 1$ , one has the Bernstein condition*

$$\sup_{g \in \mathcal{G}} P|g|^m \leq \frac{m!}{2}K^{m-2}, \quad m = 2, 3, \dots \quad (11)$$

Let  $S$  be some integer,  $\tau := 3\sqrt{6}$  and  $\delta := 4\sqrt{n} \sum_{s=1}^S 2^{-2s}/K_{s-1} + \sqrt{n}2^{-S}$ , where  $\{K_{s-1}\}_{s=1}^S$  is an arbitrary decreasing sequence of positive constants (called truncation levels). Suppose that  $\tilde{N}_s < \infty$  for all  $s = 0, \dots, S$ . Then there is a  $(\delta, \tau, \mathcal{L})$  finite tree chain for  $\{\nu_n(g)\}$ , with  $|G_s| \leq N_s$ ,  $s = 0, \dots, S$ , and with

$$L_0 = \frac{4\sqrt{6}K}{\sqrt{n}}, \quad L_s = \frac{2\sqrt{6} 2^s K_{s-1}}{3\sqrt{n}}, \quad s = 1, \dots, S.$$

As a consequence, we can derive a bound for the expectation of the supremum of the empirical process.

**Theorem 7** *Assume the Bernstein condition (11). Let*

$$\bar{\mathbf{E}}_S := 2^{-S}\sqrt{n} + 14 \sum_{s=0}^S 2^{-s} \sqrt{6\tilde{H}_s} + 6^2 K \frac{\tilde{H}_0}{\sqrt{n}}.$$

Then one has

$$\mathbb{E} \left( \sup_{g \in \mathcal{G}} |\nu_n(g)| \right) \leq \min_S \bar{\mathbf{E}}_S.$$

*Remark 3* When  $\Theta$  is finite, say  $|\Theta| = p$ , one may choose a bound with  $S = \delta = 0$ , and  $\tilde{H}_0 \leq \log(1 + p)$ . Theorem then 7 yields - up to constants - the same bound as in (2).

Finally, we present the main result of this section. We give deviation results in probability and in Bernstein-Orlicz norm, where the dependency on the complexity of  $\mathcal{G}$  is only in the shift.

**Theorem 8** *Assume the Bernstein condition (11). Define as in Theorem 7,*

$$\bar{\mathbf{E}}_S := 2^{-S}\sqrt{n} + 14 \sum_{s=0}^S 2^{-s} \sqrt{6\tilde{H}_s} + 6^2 K \frac{H_0}{\sqrt{n}}.$$

Let

$$\tilde{L} := \frac{\sqrt{6}K}{2\sqrt{n}}.$$

Then for all  $t > 0$ ,

$$\mathbf{P}\left(\sup_{g \in \mathcal{G}} |\nu_n(g)| \geq \min_S \bar{\mathbf{E}}_S + 6^2 K/\sqrt{n} + 24\sqrt{6}\left[\sqrt{t} + \frac{\tilde{L}t}{2}\right]\right) \leq 2 \exp[-t].$$

Moreover,

$$\left\| \left( \left[ \sup_{g \in \mathcal{G}} |\nu_n(g)| \right] - \left[ \min_S \bar{\mathbf{E}}_S + 6^2 K/\sqrt{n} + 24\sqrt{6} \right] \right)_+ \right\|_{\Psi_{\sqrt{3}\tilde{L}}} \leq 72\sqrt{2}.$$

Theorem 8 can be compared to results in Adamczak [2008]. One sees that our bound replaces the sub-exponential Orlicz-norm

$$\left\| \max_{1 \leq i \leq n} \sup_{g \in \mathcal{G}} |g(X_i)| \right\|_{\Psi}, \quad \Psi(z) = \exp(z) - 1, z \geq 0,$$

occurring in Adamczak [2008] by a constant proportional to  $K$ , which means we generally gain a  $\log n$ -term. On the other hand, the shift in Adamczak [2008] is up to a factor  $(1 + \epsilon)$  equal to the expectation

$$\mathbf{E} \sup_{g \in \mathcal{G}} |\nu_n(g)|,$$

as in Massart [2000]) (whose result is cited here in Theorem 5).

*Remark 4* Again, when  $|\Theta| = p$  is finite, one can choose  $S = \delta = 0$ , and  $\tilde{H}_0 \leq \log(1 + p)$ . as in Remark 3. Theorem 8 then reduces to the usual union bound type deviation inequalities for the maximum of finitely many random variables (that is, the results are - up to constants - a special case of Lemmas 4 and 5).

## 6 Proofs for Section 5

### 6.1 Proof of Theorem 6

This follows from similar arguments as in van de Geer [2000], who uses in turn ideas of Ossiander [1987]. Let for  $s = 1, \dots, S$ ,

$$\{[\tilde{g}_j^{s,L}, \tilde{g}_j^{s,U}]\}_{j=1}^{\tilde{N}_s}$$

be a minimal  $2^{-s}$ -bracketing set for  $\|\cdot\|$ . Let  $\{[\tilde{g}_j^{0,L}, \tilde{g}_j^{0,U}]\}_{j=1}^{\tilde{N}_0}$  be a generalized bracketing set.

Consider some  $g \in \mathcal{G}$ , and let  $[\tilde{g}^{0,L}, \tilde{g}^{0,U}]$  be the corresponding generalized bracket, and for all  $s \in \{1, \dots, S\}$ , let the corresponding brackets be  $[\tilde{g}^{s,L}, \tilde{g}^{s,U}]$ . Thus

$$\tilde{g}^{s,L} \leq g \leq \tilde{g}^{s,U}, \quad s = 0, \dots, S,$$

and

$$P|\tilde{g}_{0,U} - \tilde{g}_{0,L}|^m \leq \frac{m!}{2}(2K)^{m-2}, \quad m = 2, 3, \dots,$$

$$P|\tilde{g}^{s,U} - \tilde{g}^{s,L}|^2 \leq 2^{-2s}, \quad s = 1, \dots, S.$$

If for some  $s$  there are several brackets in  $\{[\tilde{g}_j^{s,L}, \tilde{g}_j^{s,U}]\}_{j=1}^{\tilde{N}_s}$  corresponding to  $g$ , we choose a fixed but otherwise arbitrary one. Define

$$g^{s,L} := \max_{0 \leq k \leq s} \tilde{g}^{k,L}, \quad g^{s,U} := \min_{0 \leq k \leq s} \tilde{g}^{k,U}.$$

Then

$$g^{0,L} \leq g^{1,L} \leq \dots \leq g^{S,L} \leq g \leq g^{S,U} \leq \dots \leq g^{1,U} \leq g^{0,U},$$

and moreover  $g^{s,U} - g^{s,L} \leq \tilde{g}^{s,U} - \tilde{g}^{s,L}$ . Denote the difference between upper and lower bracket by

$$\Delta^s := g^{s,U} - g^{s,L}, \quad s = 0, \dots, S.$$

The differences  $\Delta^s$  are decreasing in  $s$ . Furthermore,  $\|\Delta^s\| \leq 2^{-s}$ , for all  $s \in \{0, 1, \dots, S\}$ .

Let  $\mathcal{N}_s := |\{[g_j^{s,L}, g_j^{s,U}]\}|$ ,  $s = 0, \dots, S$ . It is easy to see that

$$\mathcal{N}_s \leq \prod_{k=0}^s \tilde{N}_k =: N_s, \quad s = 0, \dots, S.$$

We define a tree with end nodes  $\{1, \dots, \mathcal{N}_S\}$ . At each end node  $j$  sits a pair of brackets  $[g_j^{S,L}, g_j^{S,U}]$ . For each  $s = 0, \dots, S-1$ , we define the parents at generation  $s$  as follows. Let

$$\tilde{V}_k^s := \{l : [g_k^{s-1,L}, g_k^{s-1,U}] \text{ forms a } 2^{-(s-1)\text{-bracket for } [g_l^{s,L}, g_l^{s,U}]\}.$$

Then  $\cup_{k=1}^{\mathcal{N}_{s-1}} \tilde{V}_k^s = \{1, \dots, \mathcal{N}_s\}$ , that is, for each bracket  $[g_l^{s,L}, g_l^{s,U}]$  there is a  $k \in \{1, \dots, \mathcal{N}_{s-1}\}$  with  $l \in \tilde{V}_k^s$ . To see this, we note that for each  $l$ , there is a function  $g$  with  $g_l^{s,L} \leq g \leq g_l^{s,U}$ , and by the above construction, there is a  $k$  with  $g_k^{s-1,L} \leq g_l^{s,L} \leq g \leq g_l^{s,U} \leq g_k^{s-1,U}$ . We let  $\{V_k^s\}_{k=1}^{\mathcal{N}_{s-1}}$  be a disjoint version of  $\{\tilde{V}_k^s\}$ , e.g., the one given by

$$V_1^s = \tilde{V}_1^s, \quad V_k^s = \tilde{V}_k^s \setminus \cup_{l=1}^{k-1} \tilde{V}_l^s, \quad k = 1, \dots, \mathcal{N}_{s-1}.$$

We let

$$\text{parent}(j_s) = k \text{ if } j_s \in V_k^s.$$

We now turn to an adaptive truncation device. For for each  $s = 0, \dots, S-1$ , we are given truncation levels  $K_s$ , such that  $K_s$  is assumed to be decreasing in  $s$ . Let  $g$  be fixed and

$$g^{0,L} \leq g^{1,L} \leq \dots \leq g^{S,L} \leq g \leq g^{S,U} \leq \dots \leq g^{1,U} \leq g^{0,U}.$$

Define

$$\Delta^s := g^{s,U} - g^{s,L}, \quad y_s := 1\{\Delta^s \geq K_s\}.$$

Then

$$K_s 1\{y_s = 1\} \leq \Delta^s 1\{y_s = 1\}, \quad s = 0, \dots, S-1,$$

which implies (for  $s = 0, \dots, S-1$ )

$$P\Delta^s \mathbb{1}\{y_s = 1\} \leq \frac{P|\Delta^s|^2}{K_s} \leq \frac{2^{-2s}}{K_s}$$

We can write any  $g \in \mathcal{G}$  as

$$\begin{aligned} g &= \sum_{s=1}^S (g - g^{0,s}) \mathbb{1}\{y_s = 1, y_{s-1} = \dots = y_0 = 0\} \\ &+ \sum_{s=1}^S (g^{s,L} - g^{s-1,L}) \mathbb{1}\{y_{s-1} = \dots = y_0 = 0\} + g_{0,L} + (g - g^{0,L}) \mathbb{1}\{y_0 = 1\} \end{aligned} \quad (12)$$

Let

$$W_{j_0} := |\nu_n(g^{0,L})| + |\nu_n(\Delta^0)|,$$

$$W_{j_s} := |\nu_n(\Delta^s \mathbb{1}\{y_{s-1} = 0\})| + |\nu_n((g^{s,L} - g^{s-1,L}) \mathbb{1}\{y_{s-1} = 0\})|, \quad s = 1, \dots, S.$$

Then it follows from (12) that

$$|\nu_n(g)| \leq \sum_{s=0}^S |W_{j_s}| + \sqrt{n} \sum_{s=0}^S P\Delta^s \mathbb{1}\{y_s = 1\} \leq \sum_{s=0}^S |W_{j_s}| + \delta,$$

for

$$\delta = \sqrt{n} \sum_{s=1}^S \frac{4 \cdot 2^{-2s}}{K_{s-1}} + \sqrt{n} 2^{-S}.$$

Note now that

$$\begin{aligned} (P|g^{0,L}|^m)^{1/m} &\leq (P|g|^m)^{1/m} + (P|\Delta^0|^m)^{1/m} \\ &\leq \left(\frac{m!}{2} K^{m-2}\right)^{1/m} + \left(\frac{m!}{2} (2K)^{m-2}\right)^{1/m} \leq 2 \left(\frac{m!}{2} (2K)^{m-2}\right)^{1/m}, \end{aligned}$$

so

$$P|g^{0,L}|^m \leq \frac{m!}{2} (4K)^{m-2} 2^2.$$

By Corollary 1

$$\|\nu_n(g^{0,L})\|_{\psi_{L_0}} \leq 2\sqrt{6},$$

for

$$L_0 = \sqrt{6}(8K/2)/\sqrt{n} = 4\sqrt{6}/\sqrt{n},$$

where we multiplied by a factor 2 because the Bernstein condition for the centered functions holds with the above  $4K$  replaced by  $8K$ . Moreover,  $L_0 = \sqrt{6}(4K)/\sqrt{n}$ , so again by Corollary 1,

$$\|\nu_n(\Delta^0)\|_{\psi_{L_0}} \leq \sqrt{6}.$$

The triangle inequality gives

$$\left\| |\nu_n(g^{0,L})| + |\nu_n(\Delta^0)| \right\|_{\psi_{L_0}} \leq 3\sqrt{6} =: \tau.$$

Moreover, for  $s = 1, \dots, S$ ,

$$|(g^{s,L} - g^{s-1,L}) \mathbb{1}\{y_{s-1} = 0\}| \leq \Delta^{s-1} \leq K_{s-1}, \quad \|\Delta^{s-1}\| \leq 2^{-s+1},$$

and

$$\Delta^s \mathbb{1}\{y_{s-1} = 0\} \leq \Delta^{s-1} \leq K_{s-1}, \quad \|\Delta^s\| \leq 2^{-s}.$$

So, again by Corollary 1, we may take

$$L_s := \sqrt{6} 2^s \max\left(\frac{2}{3}K_{s-1}/2, \frac{2}{3}K_{s-1}\right)/\sqrt{n} = \frac{2\sqrt{6}K_{s-1}}{3\sqrt{n}}, \quad s = 1, \dots, S.$$

Then, again by the triangle inequality,

$$\left\| |\nu_n((g^{s,L} - g^{s-1,L})\mathbb{1}\{y_{s-1} = 0\})| + |\nu_n(\Delta^s \mathbb{1}\{y_{s-1} = 0\})| \right\|_{\Psi_{L_s}} \leq 3\sqrt{6} 2^{-s}.$$

□

## 6.2 Three technical lemmas

To apply the result of Theorem 6, we need three technical lemmas. First we need a bound for  $N_s := \prod_{k=0}^s \tilde{N}_k$ , or actually for  $H_s := \log(1 + N_s)$ .

**Lemma 6** *Let  $s \in \{0, \dots, S\}$ ,  $H_s := \log(1 + \prod_{k=0}^s \tilde{N}_k)$  and  $\tilde{H}_s := \log(1 + \tilde{N}_s)$ . It holds that*

$$\sum_{s=1}^S 2^{-s} \sqrt{H_s} \leq \sqrt{\tilde{H}_0} + 2 \sum_{s=1}^S 2^{-s} \sqrt{\tilde{H}_s}.$$

**Proof of Lemma 6.** We have

$$\sqrt{H_s} \leq \sum_{k=0}^s \sqrt{\tilde{H}_k},$$

so

$$\begin{aligned} \sum_{s=1}^S 2^{-s} \sqrt{H_s} &\leq \sum_{s=1}^S 2^{-s} \sqrt{\tilde{H}_0} + \sum_{s=1}^S 2^{-s} \sum_{k=1}^s \sqrt{\tilde{H}_k} \\ &\leq \sqrt{\tilde{H}_0} + \sum_{k=1}^S \sum_{s=k}^S 2^{-s} \sqrt{\tilde{H}_k} \leq \sqrt{\tilde{H}_0} + 2 \sum_{k=1}^S 2^{-k} \sqrt{\tilde{H}_k}. \end{aligned}$$

□

The next lemma inserts a special choice for the truncation levels  $\{K_s\}$ , and then establishes a bound for the expectation of the supremum of the empirical process, derived from the one of Theorem 2.

**Lemma 7** *Let  $S$  be some integer and  $\epsilon \geq 0$  be an arbitrary constant. Take*

$$K_{s-1} := 2^{-s} \sqrt{n} \left( \frac{\sqrt{6}}{3\sqrt{\log(1 + N_s)}} \wedge \frac{1}{\epsilon} \right), \quad s = 1, \dots, S,$$

where  $u \wedge v$  denotes the minimum of  $u$  and  $v$ . Define as in Theorem 6,

$$L_0 := \frac{4\sqrt{6}K}{\sqrt{n}}, \quad L_s := \frac{2\sqrt{6} 2^s K_{s-1}}{3\sqrt{n}}, \quad s = 1, \dots, S,$$

$$\delta := 4\sqrt{n} \sum_{s=1}^S 2^{-2s} / K_{s-1} + \sqrt{n} 2^{-S},$$

and  $\tau := 3\sqrt{6}$ . Let

$$\mathbf{E}_S := \tau \sum_{s=0}^S 2^{-s} \left[ \sqrt{\log(1 + N_s)} + \frac{L_s}{2} \log(1 + N_s) \right] + \delta.$$

Then

$$\mathbf{E}_S \leq \bar{\mathbf{E}}_S + 4\epsilon,$$

where

$$\bar{\mathbf{E}}_S := 2^{-S} \sqrt{n} + 14 \sum_{s=0}^S 2^{-s} \sqrt{6\tilde{H}_s} + 6^2 K \frac{H_0}{\sqrt{n}}.$$

**Proof of Lemma 7.** We have

$$\begin{aligned} \mathbf{E}_S &= \sum_{s=1}^S \frac{4 \cdot 2^{-2s} \sqrt{n}}{K_{s-1}} + 2^{-S} \sqrt{n} + \tau \sqrt{\log(1 + N_0)} + 2\sqrt{6} \tau K \frac{\log(1 + N_0)}{\sqrt{n}} \\ &\quad + \tau \sum_{s=1}^S 2^{-s} \left[ \sqrt{\log(1 + N_s)} + \frac{1}{3} \sqrt{6} \cdot 2^s K_{s-1} \frac{\log(1 + N_s)}{\sqrt{n}} \right] \\ &= \sum_{s=1}^S \frac{4 \cdot 2^{-2s} \sqrt{n}}{K_{s-1}} + 2^{-S} \sqrt{n} + 3\sqrt{6 \log(1 + N_0)} + 6^2 K \frac{\log(1 + N_0)}{\sqrt{n}} \\ &\quad + 3 \sum_{s=1}^S 2^{-s} \sqrt{6 \log(1 + N_s)} + \sum_{s=1}^S 6K_{s-1} \frac{\log(1 + N_s)}{\sqrt{n}} = I + II + III, \end{aligned}$$

where

$$I := 2^{-S} \sqrt{n} + 3\sqrt{6 \log(1 + N_0)} + 6^2 K \frac{\log(1 + N_0)}{\sqrt{n}},$$

$$II := 3 \sum_{s=1}^S 2^{-s} \sqrt{6 \log(1 + N_s)},$$

and

$$III := \sum_{s=1}^S \frac{4 \cdot 2^{-2s} \sqrt{n}}{K_{s-1}} + \sum_{s=1}^S 6K_{s-1} \frac{\log(1 + N_s)}{\sqrt{n}}.$$

Insert

$$K_{s-1} = \frac{1}{3} \sqrt{6} \cdot 2^{-s} \sqrt{\frac{n}{\log(1 + N_s)}} \wedge 2^{-s} \frac{\sqrt{n}}{\epsilon}, \quad s = 1, \dots, S.$$

Note that  $K_s$  is decreasing in  $s$ . Moreover

$$\frac{4 \cdot 2^{-2s} \sqrt{n}}{K_{s-1}} + 6K_{s-1} \frac{\log(1 + N_s)}{\sqrt{n}} \leq 4\sqrt{6} \cdot 2^{-s} \sqrt{\log(1 + N_s)} + 4 \cdot 2^{-s} \sqrt{n} \epsilon.$$

We find

$$III \leq 4\sqrt{6} \sum_{s=1}^S 2^{-s} \sqrt{\log(1 + N_s)} + 4\epsilon,$$

so that

$$II + III \leq 7\sqrt{6} \sum_{s=1}^S 2^{-s} \sqrt{\log(1 + N_s)} + 4\epsilon.$$

Now apply Lemma 6. This gives

$$II + III \leq 7\sqrt{6}\sqrt{\log(1 + \tilde{N}_0)} + 14\sqrt{6}\sum_{s=1}^S 2^{-s}\sqrt{\log(1 + \tilde{N}_s)} + 4\epsilon.$$

Hence,

$$\begin{aligned} I + II + III &\leq 2^{-S}\sqrt{n} + 6^2K\frac{\log(1 + \tilde{N}_0)}{\sqrt{n}} + 10\sqrt{6}\sqrt{\log(1 + \tilde{N}_0)} \\ &\quad + 14\sqrt{6}\sum_{s=1}^S 2^{-s}\sqrt{\log(1 + \tilde{N}_s)} + 4\epsilon \\ &\leq 2^{-S}\sqrt{n} + 14\sqrt{6}\sum_{s=0}^S 2^{-s}\sqrt{\log(1 + \tilde{N}_s)} + 6^2K\frac{\log(1 + \tilde{N}_0)}{\sqrt{n}} + 4\epsilon. \end{aligned}$$

□

We now derive some bounds which will be used for obtaining the deviation inequalities in probability and in Bernstein-Orlicz norm of Theorem 8.

**Lemma 8** *Let the constants  $\{K_{s-1}\}_{s=1}^S$ ,  $\{L_s\}_{s=0}^S$ , and  $\tau$  be as in Lemma 8. Let*

$$L := \sum_{s=0}^S 2^{-s} \frac{L_s(1+s)}{4}.$$

Then

$$L \leq \sqrt{6}K/\sqrt{n} + 2 \wedge \frac{\sqrt{6}}{\epsilon},$$

and

$$4\tau(1 + L/2) \leq 6^2K/\sqrt{n} + 24\sqrt{6}.$$

**Proof of Lemma 8 .** We have

$$\begin{aligned} L &= \frac{L_0}{4} + \sum_{s=1}^S \frac{2^{-s}L_s(1+s)}{4} \\ &= \frac{\sqrt{6}K}{\sqrt{n}} + \sum_{s=1}^S \frac{(1+s)K_{s-1}}{\sqrt{6n}}. \end{aligned}$$

But

$$\sum_{s=1}^S 2^{-s}(1+s) \leq 2 \int_0^\infty 2^{-x}x dx = \frac{2}{(\log 2)^2},$$

and since  $H_s = \log(1 + N_s) \geq \log(2)$ ,

$$K_{s-1} \leq 2^{-s}\sqrt{n} \left( \frac{\sqrt{6}}{3(\log(2))^{1/2}} \wedge \frac{1}{\epsilon} \right).$$

Hence,

$$\begin{aligned} L &\leq \frac{\sqrt{6}K}{\sqrt{n}} + \frac{2}{\sqrt{6}(\log 2)^2} \left( \frac{\sqrt{6}}{3(\log(2))^{1/2}} \wedge \frac{1}{\epsilon} \right) \\ &= \frac{\sqrt{6}K}{\sqrt{n}} + \frac{2}{3(\log 2)^{(5/2)}} \wedge \frac{2}{6(\log 2)^2} \frac{\sqrt{6}}{\epsilon} \\ &\leq \frac{\sqrt{6}K}{\sqrt{n}} + 2 \wedge \frac{\sqrt{6}}{\epsilon}. \end{aligned}$$

As  $\tau = 3\sqrt{6}$ , we get

$$4\tau(1 + L/2) \leq 6^2K/\sqrt{n} + 24\sqrt{6}.$$

□

## 7 Proof of Theorems 7 and 8

**Proof of Theorem 7.** This follows from Theorem 2, Theorem 6, and Lemma 7 with  $\epsilon = 0$ .  $\square$

**Proof of Theorem 8.** Let  $t > 0$  be arbitrary. Note that  $\bar{\mathbf{E}}_S$  is as in Lemma 7. Apply the bounds of Lemma 8 with  $\epsilon = 3\sqrt{t}$  for the constant  $L$  defined there. Then

$$\begin{aligned} \tau(4 + 2L) + 4\epsilon + 4\tau \left[ \sqrt{t} + \frac{Lt}{2} \right] &\leq 6^2 K/\sqrt{n} + 24\sqrt{6} + 4\epsilon + 12\sqrt{6t} + 2\tau \frac{\sqrt{6}Kt}{\sqrt{n}} + 2\tau \frac{\sqrt{6t}}{\epsilon} \\ &= 6^2 K/\sqrt{n} + 6^2 Kt/\sqrt{n} + 24\sqrt{6} + 12\sqrt{6t} + 24\sqrt{t} \\ &\leq 6^2 K/\sqrt{n} + 24\sqrt{6} + 24\sqrt{6} \left[ \sqrt{t} + \frac{\tilde{L}t}{2} \right], \end{aligned}$$

where

$$\tilde{L} := \frac{\sqrt{6}K}{2\sqrt{n}}.$$

Then by Theorem 4,

$$\mathbb{P} \left( \sup_{g \in \mathcal{G}} |\nu_n(g)| \geq \min_S \bar{\mathbf{E}}_S + 6^2 K/\sqrt{n} + 24\sqrt{6} + 24\sqrt{6} \left[ \sqrt{t} + \frac{\tilde{L}t}{2} \right] \right) \leq 2 \exp[-t].$$

and by Lemma 2

$$\left\| \left( \left[ \sup_{g \in \mathcal{G}} |\nu_n(g)| \right] - \left[ \min_S \bar{\mathbf{E}}_S + 6^2 K/\sqrt{n} + 24\sqrt{6} \right] \right)_+ \right\|_{\Psi_{\sqrt{3\tilde{L}}}} \leq 72\sqrt{2}.$$

$\square$

## References

- A. Adamczak. A tail inequality for suprema of unbounded empirical processes with applications to Markov chains. *Electronic Journal of Probability*, 13(34):1000–1034, 2008.
- G. Bennet. Probability inequalities for sums of independent random variables. *Journal of the American Statistical Association*, 57:33–45, 1962.
- J.R. Blum. On the Convergence of Empiric Distribution Functions. *The Annals of Mathematical Statistics*, 26(3):527–529, 1955.
- O. Bousquet. A Bennet concentration inequality and its application to suprema of empirical processes. *Comptes Rendus de l'Académie des Sciences, Paris*, 334: 495–550, 2002.
- P. Bühlmann and S. van de Geer. *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Springer, 2011.
- R.M. Dudley. The sizes of compact subsets of hilbert space and continuity of gaussian processes. *Journal of Functional Analysis*, pages 290–330, 1967.
- M.A. Krasnosel'skii and Y.B. Rutickii. *Convex functions and Orlicz spaces*. Noordhoff Groningen, 1961.
- J. Lederer and S. van de Geer. New concentration inequalities for suprema of empirical processes, 2011. preprint.

- 
- P. Massart. About the constants in Talagrand's concentration inequalities for empirical processes. *Annals of Probability*, 28:863–884, 2000.
- M. Ossiander. A central limit theorem under metric entropy with  $L_2$  bracketing. *Annals of Probability*, 15(3):897–919, 1987.
- M. Talagrand. New concentration inequalities in product spaces. *Inventiones Mathematicae*, 126(3):505–563, 1996.
- M. Talagrand. *The generic chaining: upper and lower bounds of stochastic processes*. Springer Verlag, 2005.
- S. van de Geer. *Empirical Processes in M-Estimation*. Cambridge University Press, 2000.
- A. W. van der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes*. Springer Series in Statistics. Springer-Verlag, New York, 1996. ISBN 0-387-94640-3.
- F.G. Viens and A.B. Vizcarra. Supremum concentration inequality and modulus of continuity for sub-nth chaos processes. *Journal of Functional Analysis*, 248:1–26, 2007.