

# Combinatorial and Structural Results for $\gamma$ - $\Psi$ -dimensions

Yann Guermeur

LORIA-CNRS

Campus Scientifique, BP 239

54506 Vandœuvre-lès-Nancy Cedex, France

(e-mail: [Yann.Guermeur@loria.fr](mailto:Yann.Guermeur@loria.fr))

May 23, 2022

**Running Title:** Combinatorial and Structural Results for  $\gamma$ - $\Psi$ -dimensions

**Keywords:** margin multi-category classifiers, guaranteed risks, scale-sensitive combinatorial dimensions,  $\gamma$ - $\Psi$ -dimensions

**Mathematics Subject Classification:** 68Q32, 62H30

## Abstract

This article deals with the generalization performance of margin multi-category classifiers, when minimal learnability hypotheses are made. In that context, the derivation of a guaranteed risk is based on the handling of capacity measures belonging to three main families: Rademacher/Gaussian complexities, metric entropies and scale-sensitive combinatorial dimensions. The usefulness of the scale-sensitive combinatorial dimensions rests on the availability of two types of results. Combinatorial results connect them to metric entropies. Structural results perform the transition from the multi-class case to the bi-class one. The results currently available for the standard scale-sensitive combinatorial dimension, the fat-shattering dimension, make it useless. We establish the advantages springing from replacing it with two  $\gamma$ - $\Psi$ -dimensions: the margin Graph dimension and the margin Natarajan dimension. Two major conclusions can be drawn:

1. involving the margin Graph dimension always improves the combinatorial results;
2. the margin Natarajan dimension can be used to exploit basic features of the classifier so as to bypass the main weakness of the fat-shattering dimension: its structural result.

## 1 Introduction

One of the main open problems of the theory of margin multi-category pattern classification is the characterization of the way the confidence interval of an upper bound on the probability of error should vary as a function of the three basic parameters which are the sample size  $m$ , the number  $C$  of categories and the margin parameter  $\gamma$  (see Kontorovich and Weiss, 2014, for a survey). When working under minimal learnability hypotheses, the derivation of such a *guaranteed risk* is based on the handling of capacity measures belonging to three main families: Rademacher/Gaussian complexities (Bartlett and Mendelson, 2002), metric entropies (Kolmogorov and Tihomirov, 1961) and scale-sensitive combinatorial dimensions (Kearns and Schapire, 1994; Guermeur, 2007). The usefulness of the scale-sensitive combinatorial dimensions to derive guaranteed risks rests on the availability of two types of results. *Combinatorial results* (Alon et al., 1997; Mendelson and Vershynin, 2003; Rudelson and Vershynin, 2006; Musayeva et al., 2019) connect them to metric entropies. *Structural results* (Duan, 2012; Maurer, 2016; Guermeur, 2017; Mohri et al., 2018) perform the transition from the multi-class case to the bi-class one. The structural result dedicated to the main scale-sensitive combinatorial dimension, the fat-shattering dimension (Kearns and

Schapire, 1994), is of no use to derive a bound. This motivates the exploration of an alternative option: replacing this dimension with the two main  $\gamma$ - $\Psi$ -dimensions (Guermeur, 2007), i.e., the margin Graph dimension and the margin Natarajan dimension. This article introduces the corresponding combinatorial and structural results. The dependence on  $m$ ,  $C$  and  $\gamma$  of the resulting guaranteed risks is then characterized. This establishes that involving the margin Graph dimension always improves the combinatorial results. Furthermore, the margin Natarajan dimension appears very promising to take into account basic features of the classifier so as to produce efficient structural results. When this happens, then the improvement of the confidence interval primarily regards the dependence on  $\gamma$ .

The organization of the paper is as follows. Section 2 introduces the theoretical framework. Section 3 highlights the need for new structural results to improve the multi-class bounds. Section 4 characterizes the connections between the three combinatorial dimensions considered. Sections 5 and 6 introduce and discuss the new combinatorial and structural results dedicated to the two  $\gamma$ - $\Psi$ -dimensions. The corresponding bounds on the metric entropies and guaranteed risks are derived in Section 7. At last, we draw conclusions in Section 8. To make reading easier, all technical lemmas and proofs have been gathered in appendix.

## 2 Margin Multi-category Classifiers

We work under minimal assumptions on the data and the classifiers, which exhibit one important feature: for each description, they return one score per category.

### 2.1 Theoretical Framework

Let  $\llbracket n_-; n_+ \rrbracket$  denote the set of integers ranging from  $n_-$  to  $n_+$ . We consider the case of  $C$ -category pattern classification problems with  $C \in \mathbb{N} \setminus \llbracket 0; 2 \rrbracket$ .  $\mathcal{X}$  is the description space and  $\mathcal{Y} = \llbracket 1; C \rrbracket$  the set of categories. Their connection is utterly characterized by an unknown probability measure  $P$ . Let  $Z = (X, Y)$  be a random pair with values in  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ , distributed according to  $P$ . We are given an  $m$ -sample  $\mathbf{Z}_m = (Z_i)_{1 \leq i \leq m} = ((X_i, Y_i))_{1 \leq i \leq m}$  made up of independent copies of  $Z$  (in short  $\mathbf{Z}_m \sim P^m$ ). The classifiers are based on classes of vector-valued functions with one component function per category. We add a basic learnability hypothesis: the classes of component functions are *uniform Glivenko-Cantelli* (uGC) (Dudley et al., 1991). Those classes must be uniformly bounded up to additive constants. We replace this property by a slightly stronger one: the vector-valued

functions take their values in a hypercube of  $\mathbb{R}^C$ . To sum up, we make minimal hypotheses to ensure that all capacity measures met in the sequel are finite (none of the bounds formulated is trivial).

**Definition 1** (Margin classifier). *Let  $\{\mathcal{G}_k : 1 \leq k \leq C\}$  be a set of classes of functions from  $\mathcal{X}$  into  $[-M_G, M_G]$  with  $M_G \in [1, +\infty)$  ( $\forall k \in \llbracket 1; C \rrbracket$ ,  $\mathcal{G}_k \subset [-M_G, M_G]^\mathcal{X}$ ). These classes are supposed to be uGC. Let  $\mathcal{G}$  be a subset of  $\prod_{k=1}^C \mathcal{G}_k$  ( $\mathcal{G}$  is a class of functions from  $\mathcal{X}$  into  $[-M_G, M_G]^C$ ). An operator  $dr$  named decision rule maps every function  $g = (g_k)_{1 \leq k \leq C} \in \mathcal{G}$  to a margin multi-category classifier  $dr_g$  in  $(\mathcal{Y} \cup \{*\})^\mathcal{X}$ . For every pair  $(g, x) \in \mathcal{G} \times \mathcal{X}$ ,  $dr_g(x)$  is either the index of the component function of  $g$  taking the highest value at  $x$ , or the dummy category  $*$  in case of ex æquo.*

The generalization capabilities of such classifiers can be characterized by means of the values taken by the differences of the component functions. This calls for the definition of standard concepts of the theory of margin multi-category pattern classification.

**Definition 2** (Margin operator  $\rho$ ). *Let  $\mathcal{G}$  be a function class defined as in Definition 1. Define  $\rho$  as an operator on  $\mathcal{G}$  such that:*

$$\begin{aligned} \rho : \mathcal{G} &\longrightarrow \rho\mathcal{G} \\ g &\longmapsto \rho_g \end{aligned}$$

$$\forall (x, k) \in \mathcal{Z}, \quad \rho_g(x, k) = \frac{1}{2} \left( g_k(x) - \max_{l \neq k} g_l(x) \right).$$

*The function  $\rho_g$  is the margin function associated with  $g$ .*

The risk of  $g \in \mathcal{G}$  is given by:  $L(g) = \mathbb{E}_{(X,Y) \sim P} [\mathbf{1}_{\{\rho_g(X,Y) \leq 0\}}] = P(\text{dr}_g(X) \neq Y)$ .

**Definition 3** (Margin loss functions). *A class of margin loss functions  $\phi_\gamma$  parameterized by  $\gamma \in (0, 1]$  is a class of nonincreasing functions from  $\mathbb{R}$  into  $[0, 1]$  satisfying:*

$$\begin{cases} \forall \gamma \in (0, 1], \phi_\gamma(0) = 1 \text{ and } \phi_\gamma(\gamma) = 0 \\ \forall (\gamma, \gamma') \in (0, 1]^2, \gamma < \gamma' \implies \phi_{\gamma'} \text{ majorizes } \phi_\gamma \end{cases}.$$

Given  $\phi_\gamma$ , the risk with margin  $\gamma$  of  $g$ ,  $L_\gamma(g)$ , is defined as:  $L_\gamma(g) = \mathbb{E}_{Z \sim P} [\phi_\gamma \circ \rho_g(Z)]$ .  $L_{\gamma,m}(g)$  designates the corresponding empirical risk, measured on  $\mathbf{Z}_m$ . When using  $\phi_\gamma$ , the behaviour of the margin functions outside the interval  $[0, \gamma]$  is irrelevant to characterize the generalization performance. The idea to exploit this property by means of a squashing operator can be traced back to Bartlett (1998). The present study uses the operator  $\pi_\gamma$ .

**Definition 4** (Squashing operator  $\pi_\gamma$ ). *Let  $\mathcal{F} \subset \mathbb{R}^{\mathcal{T}}$ . For  $\gamma \in (0, 1]$ , define the piecewise-linear squashing operator  $\pi_\gamma$  as:*

$$\begin{aligned} \pi_\gamma : \mathcal{F} &\longrightarrow \mathcal{F}_\gamma \\ f &\longmapsto f_\gamma \end{aligned}$$

$$\forall t \in \mathbb{R}, \quad f_\gamma(t) = f(t) \mathbf{1}_{\{f(t) \in (0, \gamma]\}} + \gamma \mathbf{1}_{\{f(t) > \gamma\}}.$$

When deriving a guaranteed risk, replacing the class of margin functions  $\rho_{\mathcal{G}}$  with its image by  $\pi_\gamma$ , the class  $\rho_{\mathcal{G}, \gamma}$ , has only advantages. On the one hand, it induces a decrease of the capacity (see Section 3.1) which can narrow the confidence interval. On the other hand, it does not affect the data-fit term (since  $\forall \gamma \in (0, 1], \phi_\gamma \circ \pi_\gamma = \phi_\gamma$ ). Thus, making the best of it is a major challenge.

## 2.2 Guaranteed Risks

In the theoretical framework of interest, the starting point of the derivation of a guaranteed risk is a supremum inequality taking the form:

$$P^m \left\{ \sup_{g \in \mathcal{G}} (L_*(g) - L_{\gamma, m}(g)) > F_i(m, \gamma, \delta, \text{cap}(\rho_{\mathcal{G}, \gamma})) \right\} \leq \delta, \quad (1)$$

where  $L_*$  is either  $L$  or  $L_\gamma$  and the capacity measure  $\text{cap}(\rho_{\mathcal{G}, \gamma})$  involved in the expression of the function  $F_i$  depends on the choice of  $\phi_\gamma$ . Then, the problem consists in upper bounding  $\text{cap}(\rho_{\mathcal{G}, \gamma})$  as a function of the basic parameters  $m, C$  and  $\gamma$ , so that eventually, with probability  $1 - \delta$ , the supremum of the empirical process of interest is bounded from above by a function  $F_f$  of  $m, C, \gamma$  and  $\delta$  only, i.e.,

$$\sup_{g \in \mathcal{G}} (L_*(g) - L_{\gamma, m}(g)) \leq F_f(m, C, \gamma, \delta).$$

We now introduce the three types of capacity measures considered in this study, in their order of appearance in the derivation of the bounds.

**Definition 5** (Rademacher complexity). *Let  $(\mathcal{T}, \mathcal{A}_{\mathcal{T}}, P_{\mathcal{T}})$  be a probability space and let  $T$  be a random variable distributed according to  $P_{\mathcal{T}}$ . For  $n \in \mathbb{N}^*$ , let  $\mathbf{T}_n = (T_i)_{1 \leq i \leq n}$  be an  $n$ -sample made up of independent copies of  $T$  and let  $\boldsymbol{\sigma}_n = (\sigma_i)_{1 \leq i \leq n}$  be a Rademacher sequence. Let  $\mathcal{F}$  be a class of real-valued functions with domain  $\mathcal{T}$ . The empirical Rademacher complexity of  $\mathcal{F}$  given  $\mathbf{T}_n$  is*

$$\hat{R}_n(\mathcal{F}) = \mathbb{E}_{\boldsymbol{\sigma}_n \sim \{\pm 1\}^n} \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(T_i) \mid \mathbf{T}_n \right].$$

The Rademacher complexity of  $\mathcal{F}$  is

$$R_n(\mathcal{F}) = \mathbb{E}_{\mathbf{T}_n \sim P_T^n} \left[ \hat{R}_n(\mathcal{F}) \right].$$

The classes  $\mathcal{F}$  considered here are endowed with empirical (pseudo-)metrics derived from the  $L_p$ -norms. For  $n \in \mathbb{N}^*$ , let  $\mathbf{t}_n = (t_i)_{1 \leq i \leq n} \in \mathcal{T}^n$ . Then,

$$\forall (f, f') \in \mathcal{F}^2, \begin{cases} \forall p \in [1, +\infty), d_{p, \mathbf{t}_n}(f, f') = \left( \frac{1}{n} \sum_{i=1}^n |f(t_i) - f'(t_i)|^p \right)^{\frac{1}{p}} \\ d_{\infty, \mathbf{t}_n}(f, f') = \max_{1 \leq i \leq n} |f(t_i) - f'(t_i)| \end{cases}.$$

**Definition 6** (Covering numbers and metric entropy). *Let  $\mathcal{F}$  be a class of real-valued functions on  $\mathcal{T}$  endowed with the pseudo-metric  $d_{p, \mathbf{t}_n}$ . Let  $\tilde{\mathcal{F}}$  be a totally bounded subset of  $(\mathcal{F}, d_{p, \mathbf{t}_n})$ . Then for every  $\epsilon \in \mathbb{R}_+^*$ , the  $\epsilon$ -covering number  $\mathcal{N}(\epsilon, \tilde{\mathcal{F}}, d_{p, \mathbf{t}_n})$  of  $\tilde{\mathcal{F}}$  is the minimal cardinality of a subset  $\bar{\mathcal{F}}$  of  $\tilde{\mathcal{F}}$  satisfying:*

$$\forall \tilde{f} \in \tilde{\mathcal{F}}, \exists \bar{f} \in \bar{\mathcal{F}} : d_{p, \mathbf{t}_n}(\tilde{f}, \bar{f}) < \epsilon.$$

The definition of the proper/internal covering number  $\mathcal{N}^{int}(\epsilon, \tilde{\mathcal{F}}, d_{p, \mathbf{t}_n})$  results from the restriction  $\bar{\mathcal{F}} \subset \tilde{\mathcal{F}}$ . The uniform covering numbers  $\mathcal{N}_p(\epsilon, \tilde{\mathcal{F}}, n)$  and  $\mathcal{N}_p^{int}(\epsilon, \tilde{\mathcal{F}}, n)$  are given by:

$$\begin{cases} \mathcal{N}_p(\epsilon, \tilde{\mathcal{F}}, n) = \sup_{\mathbf{t}_n \in \mathcal{T}^n} \mathcal{N}(\epsilon, \tilde{\mathcal{F}}, d_{p, \mathbf{t}_n}) \\ \mathcal{N}_p^{int}(\epsilon, \tilde{\mathcal{F}}, n) = \sup_{\mathbf{t}_n \in \mathcal{T}^n} \mathcal{N}^{int}(\epsilon, \tilde{\mathcal{F}}, d_{p, \mathbf{t}_n}) \end{cases}.$$

The function mapping  $\epsilon$  to the binary logarithm of the  $\epsilon$ -covering number of a set is called the metric entropy of this set.

There is a close connection between covering and packing properties of bounded subsets in pseudo-metric spaces.

**Definition 7** (Packing numbers). *Let the pseudo-metric space  $(\mathcal{F}, d_{p, \mathbf{t}_n})$  be defined as in Definition 6. Then for every  $\epsilon \in \mathbb{R}_+^*$ , its subset  $\tilde{\mathcal{F}}$  is  $\epsilon$ -separated with respect to  $d_{p, \mathbf{t}_n}$  if and only if:*

$$\forall \{f, f'\} \subset \tilde{\mathcal{F}}, d_{p, \mathbf{t}_n}(f, f') \geq \epsilon.$$

If  $\tilde{\mathcal{F}}$  is totally bounded, then its  $\epsilon$ -packing number  $\mathcal{M}(\epsilon, \tilde{\mathcal{F}}, d_{p, \mathbf{t}_n})$  is the maximal cardinality of its  $\epsilon$ -separated subsets.  $\mathcal{M}_p(\epsilon, \tilde{\mathcal{F}}, n) = \sup_{\mathbf{t}_n \in \mathcal{T}^n} \mathcal{M}(\epsilon, \tilde{\mathcal{F}}, d_{p, \mathbf{t}_n})$  designates the corresponding uniform packing number.

The scale-sensitive combinatorial dimensions evaluated here are  $\gamma$ - $\Psi$ -dimensions, i.e., scale-sensitive extensions of the  $\Psi$ -dimensions (Ben-David et al., 1995).

**Definition 8** ( $\gamma$ - $\Psi$ -dimensions, Definition 28 in Guermeur, 2007). Let  $\mathcal{F} \subset \mathbb{R}^{\mathcal{Z}}$  be such that:

$$\forall f \in \mathcal{F}, \forall x \in \mathcal{X}, \max_{1 \leq k < l \leq C} \{f(x, k) + f(x, l)\} = 0. \quad (2)$$

Let  $\Psi$  be a family of mappings from  $\mathcal{Y}$  into  $\{-1, 0, 1\}$ . For  $\gamma \in \mathbb{R}_+^*$ , a subset  $s_{\mathcal{Z}^n} = \{z_i = (x_i, y_i) : 1 \leq i \leq n\}$  of  $\mathcal{Z}$  is said to be  $\gamma$ - $\Psi$ -shattered by  $\mathcal{F}$  if there is a vector  $\boldsymbol{\psi}_n = (\psi^{(i)})_{1 \leq i \leq n} \in \Psi^n$  satisfying  $(\psi^{(i)}(y_i))_{1 \leq i \leq n} = \mathbf{1}_n$ , and a vector  $\mathbf{b}_n = (b_i)_{1 \leq i \leq n} \in \mathbb{R}_+^n$  such that, for every vector  $\mathbf{s}_n = (s_i)_{1 \leq i \leq n} \in \{-1, 1\}^n$ , there is a function  $f_{\mathbf{s}_n} \in \mathcal{F}$  satisfying

$$\forall i \in \llbracket 1; n \rrbracket, s_i \left( s_i \max_{\{k: \psi^{(i)}(k)=s_i\}} f_{\mathbf{s}_n}(x_i, k) - b_i \right) \geq \gamma. \quad (3)$$

The  $\gamma$ - $\Psi$ -dimension of  $\mathcal{F}$ , denoted by  $\gamma$ - $\Psi$ - $\dim(\mathcal{F})$ , is the maximal cardinality of a subset of  $\mathcal{Z}$   $\gamma$ - $\Psi$ -shattered by  $\mathcal{F}$ , if such maximum exists. Otherwise,  $\mathcal{F}$  is said to have infinite  $\gamma$ - $\Psi$ -dimension.

**Remark 1.** Let us consider the degenerate case  $C = 2$ . Then,  $s_i \max_{\{k: \psi^{(i)}(k)=s_i\}} f_{\mathbf{s}_n}(x_i, k) = f_{\mathbf{s}_n}(z_i)$ , so that Formula (3) reduces to

$$\forall i \in \llbracket 1; n \rrbracket, s_i (f_{\mathbf{s}_n}(z_i) - b_i) \geq \gamma,$$

and thus Definition 8 reduces to the definition of the main scale-sensitive combinatorial dimension, the fat-shattering or  $\gamma$ -dimension  $\gamma$ - $\dim$  (Kearns and Schapire, 1994; Bartlett and Long, 1998), with a restricted domain for vector  $\mathbf{b}_n$ . Furthermore, if we define the function class  $\mathcal{F}_{(1)}$  on  $\mathcal{X}$  as follows:  $\mathcal{F}_{(1)} = \{f(\cdot, 1) : f \in \mathcal{F}\}$ , then

$$\forall \gamma \in \mathbb{R}_+^*, \gamma\text{-dim}(\mathcal{F}_{(1)}) = \gamma\text{-dim}(\mathcal{F}),$$

provided that the definition of  $\gamma$ - $\dim(\mathcal{F}_{(1)})$  is the standard one (requiring only that  $\mathbf{b}_n \in \mathbb{R}^n$ ). Thus, in the bi-class case, the constraint  $\mathbf{b}_n \in \mathbb{R}_+^n$  of Definition 8 establishes the equivalence of the two definitions of the fat-shattering dimension, for the function class on  $\mathcal{Z}$  (classifier with two outputs) and the one on  $\mathcal{X}$  (classifier with one single output).

Definition 8 and Remark 1 suggest to adopt the following convention. The definition of any scale-sensitive combinatorial dimension of a class of functions with domain  $\mathcal{Z}$  includes the restriction  $\mathbf{b}_n \in \mathbb{R}_+^n$ . On the contrary, when the domain is  $\mathcal{X}$ , then the standard hypothesis  $\mathbf{b}_n \in \mathbb{R}^n$  applies. The relevance of this choice will appear gradually (sometimes implicitly) in the sequel.

**Definition 9** (Margin Graph dimension and margin Natarajan dimension). *Let  $\mathcal{F}$  be a function class defined as in Definition 8 and let  $\gamma \in \mathbb{R}_+^*$ . The Graph dimension with margin  $\gamma$  of  $\mathcal{F}$ , denoted by  $\gamma$ -G-dim( $\mathcal{F}$ ), is the  $\gamma$ - $\Psi$ -dimension of  $\mathcal{F}$  corresponding to the following choice for  $\Psi$ :*

$$\Psi_G = \{(\psi_k : y \mapsto \mathbb{1}_{\{y=k\}} - \mathbb{1}_{\{y \neq k\}}) : k \in \mathcal{Y}\}.$$

*The Natarajan dimension with margin  $\gamma$  of  $\mathcal{F}$ , denoted by  $\gamma$ -N-dim( $\mathcal{F}$ ), is the  $\gamma$ - $\Psi$ -dimension of  $\mathcal{F}$  corresponding to the following choice for  $\Psi$ :*

$$\Psi_N = \{(\psi_{k,l} : y \mapsto \mathbb{1}_{\{y=k\}} - \mathbb{1}_{\{y=l\}}) : \{k, l\} \subset \mathcal{Y}\}.$$

**Remark 2.** *The instance of (3) associated with the margin Graph dimension is obtained by setting  $\psi_n = (\psi_{y_i})_{1 \leq i \leq n}$  so that*

$$\forall i \in \llbracket 1; n \rrbracket, \begin{cases} \text{if } s_i = 1, f_{\mathbf{s}_n}(x_i, y_i) - b_i \geq \gamma \\ \text{if } s_i = -1, \max_{k \neq y_i} f_{\mathbf{s}_n}(x_i, k) + b_i \geq \gamma \end{cases}.$$

*In the case of the Natarajan dimension with margin  $\gamma$ , choosing  $\psi_n$  is equivalent to choosing a vector  $\mathbf{c}_n = (c_i)_{1 \leq i \leq n} \in \mathcal{Y}^n$  satisfying for every  $i \in \llbracket 1; n \rrbracket$ ,  $c_i \neq y_i$ . Then,  $\psi_n$  is set equal to  $(\psi_{y_i, c_i})_{1 \leq i \leq n}$ , so that (3) becomes*

$$\forall i \in \llbracket 1; n \rrbracket, \begin{cases} \text{if } s_i = 1, f_{\mathbf{s}_n}(x_i, y_i) - b_i \geq \gamma \\ \text{if } s_i = -1, f_{\mathbf{s}_n}(x_i, c_i) + b_i \geq \gamma \end{cases}.$$

### 2.3 Scheme of Derivation of the Guaranteed Risks

For all known instances of Formula (1), the scheme of derivation of function  $F_f$  involving the families of capacity measures considered in this study is standard. It corresponds to the directed graph depicted in Figure 1.

Here, the function class  $\mathcal{G}_0$  is equal to  $\bigcup_{k=1}^C \mathcal{G}_k$ . The value of  $m'$  is either  $m$  or  $2m$ , when the derivation of Inequality (1) involves a ghost sample (Vapnik and Chervonenkis, 1971; Pollard, 1984). When following a path from the source to the target, two types of transitions are met. A first group, the horizontal arrows, corresponds to a change of capacity measure. The standard sequence (from left to right) consists in the chaining method (Dudley, 1967; Talagrand, 2014), to connect the Rademacher complexity to covering numbers, a transition through the corresponding packing numbers, and then a combinatorial result, to switch to a combinatorial dimension. The second group, the layer of vertical arrows, is that of the structural results, performing the transition from the capacity of  $\rho_{\mathcal{G}, \gamma}$  to that of  $\mathcal{G}_0$



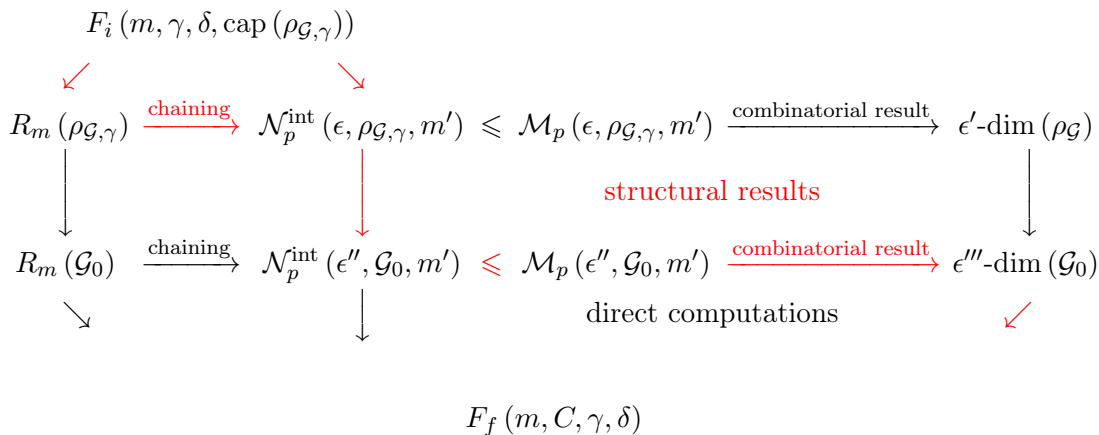


Figure 1: Graph of the transitions from a function  $F_i$  to a function  $F_f$ .

(roughly speaking from the multi-class case to the bi-class one). As an example, the paths in red are the ones explored in Guermeur (2017).

### 3 State-of-the-Art Structural Results

The literature provides us with structural results for all three types of capacity measures considered.

#### 3.1 Major Lemmas

The sharpest structural result for the Rademacher complexity of the class  $\rho_{\mathcal{G}, \gamma}$  is obtained by combining the proof of Theorem 9.2 in Mohri et al. (2018) with Talagrand's contraction lemma (see for instance Lemma 5.7 in Mohri et al., 2018).

**Lemma 1.** *Let  $\mathcal{G}$  be a function class defined as in Definition 1. Then,*

$$\forall \gamma \in (0, 1], \forall n \in \mathbb{N}^*, R_n(\rho_{\mathcal{G}, \gamma}) \leq \min \{R_n(\rho_{\mathcal{G}}), CR_n(\mathcal{G}_0)\}.$$

The counterpart of Lemma 1 dealing with covering numbers is the following structural result.

**Lemma 2** (Lemma 1 in Guermeur, 2017). *Let  $\mathcal{G}$  be a function class defined as in Definition 1. For every  $\gamma \in (0, 1]$ ,  $\epsilon \in \mathbb{R}_+^*$ ,  $n \in \mathbb{N}^*$ ,  $p \in [1, +\infty]$ , and  $\mathbf{z}_n = ((x_i, y_i))_{1 \leq i \leq n} \in \mathcal{Z}^n$ ,*

$$\mathcal{N}^{int}(\epsilon, \rho_{\mathcal{G}, \gamma}, d_{p, \mathbf{z}_n}) \leq \mathcal{N}^{int}(\epsilon, \rho_{\mathcal{G}}, d_{p, \mathbf{z}_n}) \leq \prod_{k=1}^C \mathcal{N}^{int}\left(C^{-\frac{1}{p}}\epsilon, \mathcal{G}_k, d_{p, \mathbf{x}_n}\right) \leq \left(\mathcal{N}^{int}\left(C^{-\frac{1}{p}}\epsilon, \mathcal{G}_0, d_{p, \mathbf{x}_n}\right)\right)^C,$$

where  $\mathbf{x}_n = (x_i)_{1 \leq i \leq n}$ .

The main method available to derive structural results for the  $\gamma$ -dimension (see for instance the proof of Lemma 6.2 in Duan, 2012) consists in three main steps: upper bounding the dimension of interest in terms of a metric entropy of the same class, applying a decomposition (similar to Lemma 2), and applying a combinatorial result. For the class  $\rho_{\mathcal{G}, \gamma}$ , it gives birth to the following Lemma.

**Lemma 3.** *Let  $\mathcal{G}$  be a function class defined as in Definition 1. For every  $\gamma \in (0, 1]$  and  $\epsilon \in (0, \frac{\gamma}{2}]$ ,*

$$\begin{aligned} \epsilon\text{-dim}(\rho_{\mathcal{G}, \gamma}) &\leq \epsilon\text{-dim}(\rho_{\mathcal{G}}) \\ &\leq 320 \log_2 \left( \frac{24M_{\mathcal{G}}\sqrt{C}}{\epsilon} \right) \sum_{k=1}^C \left( \frac{\epsilon}{96\sqrt{C}} \right)^k \text{-dim}(\mathcal{G}_k) \\ &\leq 320 \log_2 \left( \frac{24M_{\mathcal{G}}\sqrt{C}}{\epsilon} \right) C \left( \frac{\epsilon}{96\sqrt{C}} \right)^C \text{-dim}(\mathcal{G}_0). \end{aligned} \quad (4)$$

### 3.2 Shortcomings of the State-of-the-Art Structural Results

We reviewed the state-of-the-art decomposition results associated with the three families of capacity measures involved in this study. None is utterly satisfactory. The decomposition involving Rademacher complexities (Lemma 1) can produce a function  $F_f$  depending at least linearly on  $C$ . This behaviour is to be compared with the one of the decomposition involving covering numbers (Lemma 2), that always ensures a sublinear dependence (see for instance Theorem 3 in Musayeva et al., 2019). Furthermore, the upper bounds on the metric entropies resulting from applying a combinatorial result to  $\rho_{\mathcal{G}, \gamma}$  followed by Lemma 3 are always worse than those obtained by application of Lemma 2 followed by the same combinatorial result (applied to  $\mathcal{G}_0$ ). The reason is simple: the two computations are similar, except for an extra application of a combinatorial result in the first case. This supplementary step introduces a multiplicative factor  $\ln\left(\frac{1}{\epsilon}\right)$  in the bounds. We illustrate the phenomenon with the two state-of-the-art combinatorial results connecting the packing numbers of a class of real-valued functions to its fat-shattering dimension: Lemma 3.5 in Alon et al. (1997), for the  $L_{\infty}$ -norm, and Theorem 1 in Mendelson and Vershynin (2003), for the  $L_2$ -norm. The use of Lemma 2 produces for  $n \geq \left(\frac{\epsilon}{4}\right)\text{-dim}(\mathcal{G}_0)$ ,

$$\log_2(\mathcal{N}_{\infty}^{\text{int}}(\epsilon, \rho_{\mathcal{G}, \gamma}, n)) \leq C \left\{ \left\lceil \left(\frac{\epsilon}{4}\right)\text{-dim}(\mathcal{G}_0) \log_2 \left( \frac{4M_{\mathcal{G}}n}{\epsilon} \right) \right\rceil \log_2 \left( \frac{16M_{\mathcal{G}}^2 n}{\epsilon^2} \right) + 1 \right\} \quad (5)$$

and for  $n \in \mathbb{N}^*$ ,

$$\log_2 (\mathcal{N}_2^{\text{int}} (\epsilon, \rho_{\mathcal{G}, \gamma}, n)) \leq 20C \left( \frac{\epsilon}{48\sqrt{C}} \right) \text{-dim} (\mathcal{G}_0) \log_2 \left( \frac{12M_{\mathcal{G}}\sqrt{C}}{\epsilon} \right). \quad (6)$$

Using instead Lemma 3 gives for  $n \geq \left(\frac{\epsilon}{4}\right) \text{-dim} (\rho_{\mathcal{G}})$ ,

$$\begin{aligned} & \log_2 (\mathcal{N}_{\infty}^{\text{int}} (\epsilon, \rho_{\mathcal{G}, \gamma}, n)) \\ & \leq \left[ 320 \log_2 \left( \frac{96M_{\mathcal{G}}\sqrt{C}}{\epsilon} \right) C \left( \frac{\epsilon}{384\sqrt{C}} \right) \text{-dim} (\mathcal{G}_0) \log_2 \left( \frac{2\gamma en}{\epsilon} \right) \right] \log_2 \left( \frac{4\gamma^2 n}{\epsilon^2} \right) + 1 \end{aligned} \quad (7)$$

and for  $n \in \mathbb{N}^*$ ,

$$\log_2 (\mathcal{N}_2^{\text{int}} (\epsilon, \rho_{\mathcal{G}, \gamma}, n)) \leq 6400 \log_2 \left( \frac{1152M_{\mathcal{G}}\sqrt{C}}{\epsilon} \right) C \left( \frac{\epsilon}{4608\sqrt{C}} \right) \text{-dim} (\mathcal{G}_0) \log_2 \left( \frac{6\gamma}{\epsilon} \right). \quad (8)$$

A comparison of Inequalities (5) and (7) on the one hand, and Inequalities (6) and (8) on the other, makes it possible to identify the extra logarithmic factors. With these observations at hand, one could think that when the assumptions regarding the classifier are minimal, then the best structural result is Lemma 2. Put in a different way, the paths in the graph of Figure 1 generating the best functions  $F_f$  could be those marked in red. However, Lemma 2 makes no use of the capacity reduction induced by the squashing (operator  $\pi_{\gamma}$ ). When delaying the decomposition at this level, this squashing is only exploited upstream, i.e., by the chaining formula. Those limitations raise a question.

**Question 1.** *Can the introduction in the graph of transitions of  $\gamma$ - $\Psi$ -dimensions of  $\rho_{\mathcal{G}}$  improve the dependence of function  $F_f$  on the basic parameters?*

The answers to Question 1 should spring from replacing the classical graph of transitions (Figure 1) with the one of Figure 2 and exploring in the new graph the paths highlighted in blue.

This exploration takes the form of the derivation of new combinatorial and structural results. The evaluation of this contribution rests on connections between the scale-sensitive combinatorial dimensions which are established in the following section.

## 4 Connections Between the Scale-Sensitive Combinatorial Dimensions

We first consider the connections between the fat-shattering dimension and the margin Graph dimension. It takes the form of a kind of “equivalence”.

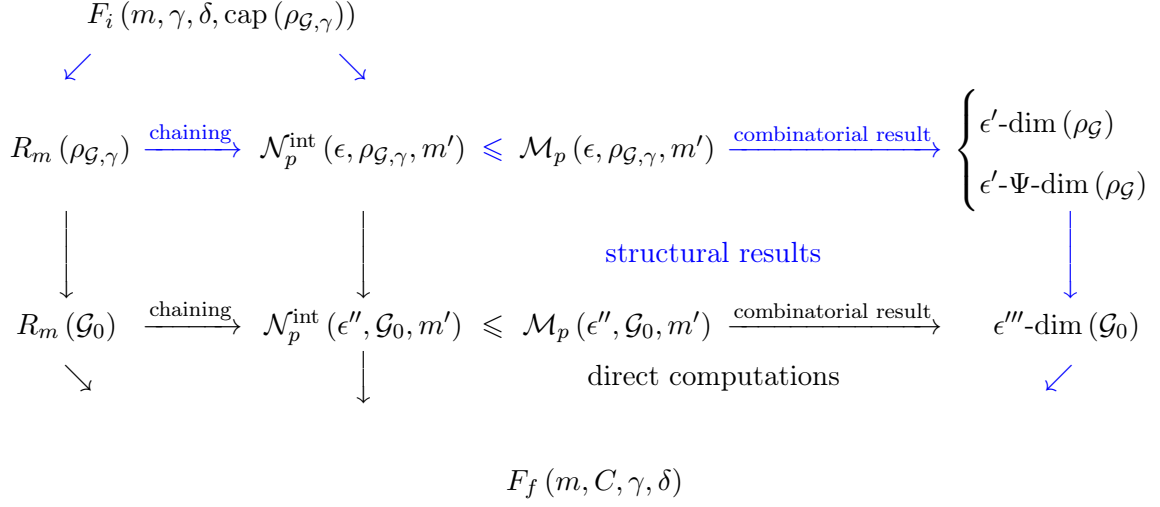


Figure 2: Paths from  $F_i$  to  $F_f$  involving combinatorial dimensions of the class  $\rho_{\mathcal{G}}$ .

**Lemma 4.** *Let  $\mathcal{F}$  be a function class defined as in Definition 8. Then,*

$$\forall \gamma \in (0, 1], \forall \epsilon \in \left(0, \frac{\gamma}{2}\right], \epsilon\text{-dim}(\mathcal{F}_\gamma) \leq \epsilon\text{-}G\text{-dim}(\mathcal{F}) \leq \epsilon\text{-dim}(\mathcal{F}). \quad (9)$$

To assess the scope of Lemma 4, one must consider the strategy implemented to establish combinatorial results for function classes  $\mathcal{F}_\gamma$  under the sole hypothesis that the fat-shattering dimension of  $\mathcal{F}$  is defined. This strategy, used for instance in Bartlett (1998), simply consists in upper bounding the packing numbers of  $\mathcal{F}_\gamma$  as a function of its fat-shattering dimension, and then upper bounding this dimension by the fat-shattering dimension of  $\mathcal{F}$ . Lemma 4 tells us that if the class  $\mathcal{F}$  also has  $\gamma$ - $\Psi$ -dimensions, then its fat-shattering dimension can be replaced with its (smaller) margin Graph dimension, thus improving the combinatorial result (whatever the  $L_p$ -norm used). The gain is conditioned by the availability of non trivial upper bounds on the margin Graph dimension of  $\mathcal{F}$ , i.e., bounds that no longer hold true if this dimension is replaced with the fat-shattering dimension of  $\mathcal{F}$ . Lemma 5, a scale-sensitive counterpart of Theorem 10 in Ben-David et al. (1995), provides a bound of this kind.

**Lemma 5.** *Let  $\mathcal{F}$  be a function class defined as in Definition 8. Suppose that  $\gamma \in \mathbb{R}_+^*$  is such that  $\gamma$ - $G$ -dim( $\mathcal{F}$ ) is finite. Then,*

$$\gamma\text{-}N\text{-dim}(\mathcal{F}) \leq \gamma\text{-}G\text{-dim}(\mathcal{F}) \leq 32 \log_2^2(e(C-1)) \gamma\text{-}N\text{-dim}(\mathcal{F})^{\alpha(C)}, \quad (10)$$

where  $\alpha(C) = 1 + \frac{1}{4 \ln(C-1) + 2}$ .

It is easy to exhibit examples of function classes  $\mathcal{F}$  with  $\gamma$ - $\Psi$ -dimensions, and for which the right-hand side inequality of Formula (10) no longer holds true with  $\gamma$ -G-dim( $\mathcal{F}$ ) replaced with  $\gamma$ -dim( $\mathcal{F}$ ). Example 1 is of this kind.

**Example 1.** Let  $\mathcal{G}$  be a set of two functions  $g^{(1)}$  and  $g^{(2)}$  on  $\mathcal{X} = \{x\}$  given by  $g^{(1)}(x) = (\frac{3}{4}, \frac{1}{4}, 0)^T$  and  $g^{(2)}(x) = (0, \frac{1}{2}, \frac{1}{2})^T$ . Let us set  $\gamma = \frac{1}{4}$ . Then,  $\gamma$ -G-dim( $\rho_{\mathcal{G}}$ ) = 0 and  $\gamma$ -dim( $\rho_{\mathcal{G}}$ ) = 1.

Indeed, by definition,  $(\rho_{g^{(1)}}(x, k))_{1 \leq k \leq 3} = (\frac{1}{4}, -\frac{1}{4}, -\frac{3}{8})^T$ , and  $(\rho_{g^{(2)}}(x, k))_{1 \leq k \leq 3} = (-\frac{1}{4}, 0, 0)^T$ . Consequently, the formula  $\max_{k \neq l} \{\rho_{g^{(1)}}(x, k) + \rho_{g^{(2)}}(x, l)\} = \gamma < 2\gamma$  implies that none of the three singletons  $\{(x, k)\}$  is  $\gamma$ -G-shattered by  $\rho_{\mathcal{G}}$ . On the contrary,

$$\begin{cases} \rho_{g^{(1)}}(x, 1) \geq \gamma \\ -\rho_{g^{(2)}}(x, 1) \geq \gamma \end{cases},$$

i.e., the class  $\rho_{\mathcal{G}}$   $\gamma$ -shatters  $\{(x, 1)\}$  for  $b = 0$ . To sum up,  $\gamma$ -dim( $\rho_{\mathcal{G}}$ ) = 1 but  $\gamma$ -N-dim( $\rho_{\mathcal{G}}$ ) =  $\gamma$ -G-dim( $\rho_{\mathcal{G}}$ ) = 0 (so that the last term of Formula (10) is also equal to 0).

## 5 Combinatorial Results

The new results exposed in this section and the next one provide the building blocks needed to derive upper bounds on the metric entropies of  $\rho_{\mathcal{G}, \gamma}$ , for  $p \geq 2$ , following the blue paths of Figure 2. However, the combinatorial results are more general, since they apply to any pair  $(\mathcal{F}_{\gamma}, \mathcal{F})$  where  $\mathcal{F}$  is a function class for which the  $\gamma$ - $\Psi$ -dimensions are defined, and not only the pairs  $(\rho_{\mathcal{G}, \gamma}, \rho_{\mathcal{G}})$ . To keep the comparison with the literature simple, we focus on the two most popular options:  $p = \infty$  and  $p = 2$ , but the generalization is straightforward using the ideas developed in the proof of Theorem 2 in Musayeva et al. (2019).

### 5.1 Margin Graph Dimension

In the case of the margin Graph dimension, the availability of Lemma 4 has a major consequence: the new combinatorial results should be compared to those of the literature applied to  $\mathcal{F}_{\gamma}$ , with the fat-shattering dimension of  $\mathcal{F}_{\gamma}$  replaced with the margin Graph dimension of  $\mathcal{F}$  (instead of its fat-shattering dimension, as is usually done). Our first result deals with the case  $p = \infty$ .

**Lemma 6.** *Let  $\mathcal{F}$  be a function class defined as in Definition 8. For  $\epsilon \in \mathbb{R}_+^*$ , let  $d_G(\epsilon) = \epsilon$ - $G$ - $\dim(\mathcal{F})$ . Then for every  $\gamma \in (0, 1]$ ,  $\epsilon \in (0, \gamma]$  and  $n \in \mathbb{N}^*$  such that  $n \geq d_G(\frac{\epsilon}{4})$ ,*

$$\mathcal{M}_\infty(\epsilon, \mathcal{F}_\gamma, n) \leq \left(\frac{6\gamma n}{\epsilon}\right)^{d_G(\frac{\epsilon}{4}) \log_2\left(\frac{2\gamma en}{d_G(\frac{\epsilon}{4})\epsilon}\right)}. \quad (11)$$

Inequality (11) compares with the application of Lemma 3.5 in Alon et al. (1997). This application produces

$$\mathcal{M}_\infty(\epsilon, \mathcal{F}_\gamma, n) < 2 \left(\frac{4\gamma^2 n}{\epsilon^2}\right)^{\left\lceil d_G(\frac{\epsilon}{4}) \log_2\left(\frac{2\gamma en}{d_G(\frac{\epsilon}{4})\epsilon}\right) \right\rceil}.$$

A gain can be noticed, which appears especially clearly for  $\epsilon = \frac{\gamma}{2}$ , the combination of practical interest as will be seen in Section 7.2. We now turn to the case  $p = 2$ .

**Lemma 7.** *Let  $\mathcal{F}$  be a function class defined as in Definition 8. For  $\epsilon \in \mathbb{R}_+^*$ , let  $d_G(\epsilon) = \epsilon$ - $G$ - $\dim(\mathcal{F})$ . Then for every  $\gamma \in (0, 1]$ ,  $\epsilon \in (0, \gamma]$  and  $n \in \mathbb{N}^*$ ,*

$$\mathcal{M}_2(\epsilon, \mathcal{F}_\gamma, n) \leq \left(\frac{5\gamma}{\epsilon}\right)^{12d_G(\frac{\epsilon}{24})}. \quad (12)$$

Inequality (12) compares with the formula obtained with Theorem 1 in Mendelson and Vershynin (2003). This time, the improvement is limited to an optimization of the constants which is not induced by the direct connection between the packing numbers of  $\mathcal{F}_\gamma$  and the margin Graph dimension of  $\mathcal{F}$ .

## 5.2 Margin Natarajan Dimension

As for the margin Natarajan dimension, with Lemma 5 at hand, Lemmas 6 and 7 also provide us with combinatorial results involving this capacity measure. However, sharper bounds should spring from following the direct path, i.e., working directly with this latter dimension (without involving the margin Graph dimension). We now state the corresponding combinatorial results (for  $p = \infty$  then  $p = 2$ ) and perform the comparison.

**Lemma 8.** *Let  $\mathcal{F}$  be a function class defined as in Definition 8. For  $\epsilon \in \mathbb{R}_+^*$ , let  $d_N(\epsilon) = \epsilon$ - $N$ - $\dim(\mathcal{F})$ . Then for every  $\gamma \in (0, 1]$ ,  $\epsilon \in (0, \gamma]$  and  $n \in \mathbb{N}^*$  such that  $n \geq d_N(\frac{\epsilon}{4})$ ,*

$$\mathcal{M}_\infty(\epsilon, \mathcal{F}_\gamma, n) \leq \left(\frac{6\gamma\sqrt{C-1}n}{\epsilon}\right)^{d_N(\frac{\epsilon}{4}) \log_2\left(\frac{2\gamma(C-1)en}{d_N(\frac{\epsilon}{4})\epsilon}\right)}. \quad (13)$$

**Lemma 9.** *Let  $\mathcal{F}$  be a function class defined as in Definition 8. For  $\epsilon \in \mathbb{R}_+^*$ , let  $d_N(\epsilon) = \epsilon$ - $N$ - $\dim(\mathcal{F})$ . Then for every  $\gamma \in (0, 1]$ ,  $\epsilon \in (0, \gamma]$  and  $n \in \mathbb{N}^*$ ,*

$$\mathcal{M}_2(\epsilon, \mathcal{F}_\gamma, n) \leq \left((C-1) \left(\frac{4\gamma}{\epsilon}\right)^5\right)^{\frac{3}{2} \log_2\left(2\left(\frac{14\gamma}{\epsilon}\right)^2(C-1)\right) d_N(\frac{\epsilon}{28})}. \quad (14)$$

As expected, as close to 1 as  $\alpha(C)$  may be, Inequalities (13) and (14) are better than the bounds obtained by substitution of the right-hand side term of (10) in the formulas involving the margin Graph dimension: (11) and (12), respectively. Precisely, in both cases, the dependence on  $n$  is unchanged, while the dependences on  $C$  and  $\epsilon$  are slightly improved.

### 5.3 Discussion

Overall, we have seen that the combinatorial results involving directly the margin Graph dimension, Lemmas 6 and 7, provide sharper bounds on the packing numbers of  $\mathcal{F}_\gamma$  (and thus also  $\rho_{\mathcal{G},\gamma}$ ) than their counterparts from the literature: Lemma 3.5 in Alon et al. (1997) and Theorem 1 in Mendelson and Vershynin (2003), respectively. The gain is more important in the first case. Furthermore, the use of the margin Natarajan dimension appears as a promising alternative. Indeed, Lemmas 8 and 9 provide sharper bounds than those resulting from the combination of Lemmas 6 and 7 with Lemma 5.

## 6 Structural Results for the Margin Natarajan Dimension

Since the margin Graph dimension is upper bounded by the fat-shattering dimension (Lemma 4), Lemma 3 is also a structural result for the margin Graph dimension. Loosely speaking, this result can be improved by substituting in its proof the  $L_2$ -norm with the norm  $L_{\lceil \log_2(C) \rceil}$ , i.e., by applying the idea developed in Musayeva et al. (2019). The gain then regards the dependence of the bound on  $C$ . However, the prohibitive drawback identified in Section 3.2 remains. It is possible to upper bound directly (without resorting to Lemma 3) the margin Graph dimension of specific classifiers from the literature, but this comes at the expense of difficult computations, that go beyond the scope of this article. In that respect, it is far easier to exploit the appealing properties of the margin Natarajan dimension, as will be seen in the sequel.

**Lemma 10.** *Let  $\mathcal{G}$  be a function class defined as in Definition 1 and let  $\mathcal{D}_{\mathcal{G}}$  be the function class  $\{\frac{1}{2}(g_k - g_l) : g \in \mathcal{G}, 1 \leq k < l \leq C\}$ . Then for every value of  $\gamma$  in  $(0, M_{\mathcal{G}}]$ ,*

$$\gamma\text{-}N\text{-dim}(\rho_{\mathcal{G}}) \leq \binom{C}{2} \cdot \gamma\text{-dim}(\mathcal{D}_{\mathcal{G}}) \quad (15)$$

and

$$\gamma\text{-}N\text{-dim}(\rho_{\mathcal{G}}) \leq 384 \binom{C}{2} \log_2 \left( \frac{20M_{\mathcal{G}}}{\gamma} \right) \left( \frac{\gamma}{48} \right)\text{-dim}(\mathcal{G}_0). \quad (16)$$

Formula (16) is almost as unsatisfactory as Formula (4). However, Formula (15) can be used to derive a sharper structural result for popular classifiers associated with function classes  $\mathcal{G}$  such that  $\mathcal{G}_0$  has the closure property  $\mathcal{D}_{\mathcal{G}} \subset \mathcal{G}_0$ . Corollary 1 illustrates this behaviour.

**Corollary 1.** *Let  $\mathcal{H}^{(1)}$  be a class of functions from  $\mathcal{X}$  into a Hilbert space  $(\mathbf{H}, \langle \cdot, \cdot \rangle_{\mathbf{H}})$  and  $(\Lambda_1, \Lambda_2) \in (\mathbb{R}_+^*)^2$ . Let  $\mathcal{H}^{(2)}$  be the class of functions  $h^{(2)}$  from  $\mathcal{X}$  into  $[-\Lambda_1\Lambda_2, \Lambda_1\Lambda_2]^C$  of the form:*

$$\forall x \in \mathcal{X}, \quad h^{(2)}(x) = \left( \left\langle \mathbf{w}_k, h^{(1)}(x) \right\rangle_{\mathbf{H}} \right)_{1 \leq k \leq C},$$

where  $h^{(1)} \in \mathcal{H}^{(1)}$  satisfies  $\sup_{x \in \mathcal{X}} \|h^{(1)}(x)\|_{\mathbf{H}} \leq \Lambda_1$  and the vector  $(\mathbf{w}_k)_{1 \leq k \leq C} \in \mathbf{H}^C$  satisfies  $\max_{1 \leq k \leq C} \|\mathbf{w}_k\|_{\mathbf{H}} \leq \Lambda_2$ . Let  $\mathcal{H}_0^{(2)}$  be the class of all the component functions of the functions in  $\mathcal{H}^{(2)}$ . Then,

$$\forall \gamma \in (0, \Lambda_1\Lambda_2], \quad \gamma\text{-N-dim}(\rho_{\mathcal{H}^{(2)}}) \leq \binom{C}{2} \cdot \gamma\text{-dim}(\mathcal{H}_0^{(2)}). \quad (17)$$

Corollary 1 is actually a consequence of Formula (15) since the class  $\mathcal{H}^{(2)}$  has been specified so as to ensure the satisfaction of the closure property. It applies to classifiers of reference such as the multi-layer perceptrons (MLPs) (Anthony and Bartlett, 1999) with linear output units and the  $C$ -category support vector machines (SVMs) (Doğan et al., 2016). In the second case,  $\mathcal{H}^{(1)}$  is restricted to one single function, the feature map, which can be defined from the kernel  $\kappa$  (Berlinet and Thomas-Agnan, 2004) as:  $\forall x \in \mathcal{X}$ ,  $h^{(1)}(x) = \kappa_x = \kappa(\cdot, x)$ . Then,  $\mathbf{H}$  is the reproducing kernel Hilbert space (RKHS) of  $\kappa$ . An application of the standard upper bound on the fat-shattering dimension of (binary) SVMs, Theorem 4.6 in Bartlett and Shawe-Taylor (1999), produces the instantiation of Inequality (17) for  $C$ -category SVMs:

$$\forall \gamma \in (0, \Lambda_1\Lambda_2], \quad \gamma\text{-N-dim}(\rho_{\mathcal{H}^{(2)}}) \leq \binom{C}{2} \left( \frac{\Lambda_1\Lambda_2}{\gamma} \right)^2. \quad (18)$$

It is noticeable that such a simple algebraic property as  $\mathcal{D}_{\mathcal{G}} \subset \mathcal{G}_0$  proves enough to replace Inequality (16) with a far sharper bound. Corollary 1 provides a first illustration of the capacity of the margin Natarajan dimension to exploit a coupling between the component functions of the classifier. This exploitation can be carried out further. We illustrate the phenomenon by refining the study of the case of the  $C$ -category SVMs. To than end, a specific definition of these machines is used, which is based on the concept of RKHS of  $\mathbb{R}^C$ -valued functions (Wahba, 1992).



**Definition 10** (RKHS  $\mathbf{H}_{\kappa,C}$ ). Let  $\kappa$  be a real-valued positive type function on  $\mathcal{X}^2$  and let  $(\mathbf{H}_{\kappa}, \langle \cdot, \cdot \rangle_{\mathbf{H}_{\kappa}})$  be its RKHS. Let  $\tilde{\kappa}$  be the real-valued positive type function on  $\mathcal{Z}^2$  deduced from  $\kappa$  as follows:  $\forall (z, z') \in \mathcal{Z}^2$ ,  $\tilde{\kappa}(z, z') = \delta_{y,y'} \kappa(x, x')$ , where  $\delta$  is the Kronecker delta. For every  $z \in \mathcal{Z}$ , let us define the  $\mathbb{R}^C$ -valued function  $\tilde{\kappa}_z^{(C)}$  on  $\mathcal{X}$  by the formula

$$\tilde{\kappa}_z^{(C)}(\cdot) = (\tilde{\kappa}(z, (\cdot, k)))_{1 \leq k \leq C}. \quad (19)$$

The RKHS of  $\mathbb{R}^C$ -valued functions at the basis of a  $C$ -category SVM with kernel  $\kappa$ ,  $(\mathbf{H}_{\kappa,C}, \langle \cdot, \cdot \rangle_{\mathbf{H}_{\kappa,C}})$ , consists of the linear manifold of all finite linear combinations of functions of the form (19) and its closure with respect to the inner product:  $\forall (z, z') \in \mathcal{Z}^2$ ,  $\langle \tilde{\kappa}_z^{(C)}, \tilde{\kappa}_{z'}^{(C)} \rangle_{\mathbf{H}_{\kappa,C}} = \tilde{\kappa}(z, z')$ .

With Definition 10 at hand, the specification of the function class at the basis of a  $C$ -category SVM rests on the condition controlling the capacity through a coupling between the component functions. We consider the standard one, used for instance by Lei et al. (2015).

**Definition 11** (Function class  $\mathcal{H}_{\Lambda}$ ). Let  $\kappa$  be a kernel on  $\mathcal{X}^2$  and let  $\Lambda \in \mathbb{R}_+^*$ . Let  $(\mathbf{H}_{\kappa,C}, \langle \cdot, \cdot \rangle_{\mathbf{H}_{\kappa,C}})$  be the RKHS of  $\mathbb{R}^C$ -valued functions spanned by  $\kappa$  according to Definition 10. Then the function class  $\mathcal{H}_{\Lambda}$  associated with the  $C$ -category SVM parameterized by  $(\kappa, \Lambda)$  is:  $\mathcal{H}_{\Lambda} = \left\{ h = (h_k)_{1 \leq k \leq C} \in \mathbf{H}_{\kappa,C} : \sum_{k=1}^C h_k = \mathbf{0}_{\mathbf{H}_{\kappa}} \text{ and } \|h\|_{\mathbf{H}_{\kappa,C}} \leq \Lambda \right\}$ .

The class  $\mathcal{H}_{\Lambda}$  can be seen as an instance of the class  $\mathcal{H}^{(2)}$  of Corollary 1 for which  $\Lambda_1 = \sup_{x \in \mathcal{X}} \|\kappa_x\|_{\mathbf{H}_{\kappa}}$  and  $\Lambda_2 = 2^{-\frac{1}{2}} \Lambda$ . This instance exhibits a stronger coupling between the component functions. Lemma 11 takes this property into account to produce a sharper bound on the margin Natarajan dimension.

**Lemma 11.** For  $\Lambda \in \mathbb{R}_+^*$ , let  $\mathcal{H}_{\Lambda}$  be the function class of Definition 11. Suppose that for every  $x \in \mathcal{X}$ ,  $\kappa_x$  belongs to the closed ball of radius  $\Lambda_{\mathcal{X}}$  about the origin in  $\mathbf{H}_{\kappa}$ . Then,

$$\forall \gamma \in (0, \Lambda \Lambda_{\mathcal{X}}], \gamma\text{-N-dim}(\rho_{\mathcal{H}_{\Lambda}}) \leq C \left( \frac{\Lambda \Lambda_{\mathcal{X}}}{2\gamma} \right)^2. \quad (20)$$

In words, the stronger coupling between the component functions could be exploited so as to turn the quadratic dependence on  $C$  of Formula (18) into a linear one.

## 7 Guaranteed Risks

The combinatorial and structural results of the two previous sections provide us with new upper bounds on the metric entropies which are based on the margin Natarajan

dimension. Their comparison with the bounds of reference, for instance Inequalities (5) and (6), requires two generic formulas. The first one is an upper bound on the  $\gamma$ -dimension of the class  $\mathcal{G}_0$ . The second one is the corresponding structural result for the margin Natarajan dimension of  $\rho_{\mathcal{G}}$ .

## 7.1 Bounds on the Metric Entropies

For the first formula, we use the standard hypothesis: that of polynomial  $\gamma$ -dimensions (van der Vaart and Wellner, 1996; Mendelson, 2003). We have already seen that it is satisfied by SVMs (Formula (18)). This is also the case for MLPs with linear output units (see for instance Theorem 14.19 in Anthony and Bartlett, 1999). The second formula is designed to incorporate the hypothesis of polynomial  $\gamma$ -dimensions in a decomposition result taking benefit from the coupling between the component functions of the functions  $g$ . It is thus primarily inspired by Inequalities (17) and (20).

**Hypothesis 1.** *We consider function classes  $\mathcal{G}$  defined as in Definition 1 for which there exists a quadruplet  $(d_{\mathcal{G},C}, d_{\mathcal{G},\gamma}, K_{\mathcal{G}_0}, K_{\rho_{\mathcal{G}}}) \in (0, 2] \times (\mathbb{R}_+^*)^3$  such that*

$$\forall \epsilon \in (0, M_{\mathcal{G}}], \begin{cases} \epsilon \text{-dim}(\mathcal{G}_0) \leq K_{\mathcal{G}_0} \epsilon^{-d_{\mathcal{G},\gamma}} & (21a) \\ \epsilon \text{-N-dim}(\rho_{\mathcal{G}}) \leq K_{\rho_{\mathcal{G}}} C^{d_{\mathcal{G},C}} \epsilon^{-d_{\mathcal{G},\gamma}}. & (21b) \end{cases}$$

Under Hypothesis 1, the combinatorial results dedicated to the margin Natarajan dimension (Lemmas 8 and 9) give birth to the following bounds on the metric entropies.

**Theorem 1.** *Let  $\mathcal{G}$  be a function class satisfying Hypothesis 1. For every  $\gamma \in (0, 1]$ ,  $\epsilon \in (0, \gamma]$  and  $n \in \mathbb{N}^*$  such that  $n \geq K_{\rho_{\mathcal{G}}} C^{d_{\mathcal{G},C}} \left(\frac{4}{\epsilon}\right)^{d_{\mathcal{G},\gamma}}$ ,*

$$\log_2(\mathcal{N}_{\infty}^{int}(\epsilon, \rho_{\mathcal{G},\gamma}, n)) \leq K_{\rho_{\mathcal{G}}} C^{d_{\mathcal{G},C}} \log_2^2\left(\frac{6\gamma(C-1)n}{\epsilon}\right) \left(\frac{4}{\epsilon}\right)^{d_{\mathcal{G},\gamma}}. \quad (22)$$

*For every  $\gamma \in (0, 1]$ ,  $\epsilon \in (0, \gamma]$  and  $n \in \mathbb{N}^*$ ,*

$$\log_2(\mathcal{N}_2^{int}(\epsilon, \rho_{\mathcal{G},\gamma}, n)) \leq \frac{3}{2} K_{\rho_{\mathcal{G}}} C^{d_{\mathcal{G},C}} \log_2^2\left((C-1) \left(\frac{4\gamma}{\epsilon}\right)^5\right) \left(\frac{28}{\epsilon}\right)^{d_{\mathcal{G},\gamma}}. \quad (23)$$

As expected, those two bounds are significantly better than the bounds obtained with the structural result dedicated to the fat-shattering dimension of  $\rho_{\mathcal{G}}$ : Lemma 3. Thus, a partial answer to Question 1 emerges: the margin Natarajan dimension is the first scale-sensitive combinatorial dimension which can be considered to handle the class of margin functions  $\rho_{\mathcal{G}}$  for a vast family of function classes  $\mathcal{G}$ . The remaining comparison to be done is with the bounds using the structural result involving covering numbers (Lemma 2), i.e.,

Inequalities (5) and (6). The rest of the section is devoted to this comparison. To make it more concrete, we use as touchstones the functions  $F_i$  corresponding to the state-of-the-art basic supremum inequalities associated with the two most popular margin loss functions.

## 7.2 Guaranteed Risk for the Indicator Margin Loss Function

The class of margin loss functions contains two main families: the family of indicator functions and the one of Lipschitz continuous functions. Our first guaranteed risk involves the classical indicator function:  $\phi_{\infty, \gamma}$ . It is given by:

$$\forall t \in \mathbb{R}, \quad \phi_{\infty, \gamma}(t) = \mathbf{1}_{\{t < \gamma\}}.$$

To the best of our knowledge, the sharpest instance of Inequality (1) involving this loss function is provided by Theorem 2 in Guermeur (2017). This bound corresponds to the case  $L_* = L$  and produces:

$$F_i(m, \gamma, \delta, \text{cap}(\rho_{\mathcal{G}, \gamma})) = \sqrt{\frac{2}{m} \left( \ln \left( \mathcal{N}_{\infty}^{\text{int}} \left( \frac{\gamma}{2}, \rho_{\mathcal{G}, \gamma}, 2m \right) \right) + \ln \left( \frac{2}{\delta} \right) \right)} + \frac{1}{m}. \quad (24)$$

By application of (22) (for  $m \geq \frac{1}{2} K_{\rho_{\mathcal{G}}} C^{d_{\mathcal{G}, C}} \left( \frac{8}{\gamma} \right)^{d_{\mathcal{G}, \gamma}}$ ):

$$\log_2 \left( \mathcal{N}_{\infty}^{\text{int}} \left( \frac{\gamma}{2}, \rho_{\mathcal{G}, \gamma}, 2m \right) \right) \leq K_{\rho_{\mathcal{G}}} C^{d_{\mathcal{G}, C}} \log_2^2 (24(C-1)m) \left( \frac{8}{\gamma} \right)^{d_{\mathcal{G}, \gamma}}. \quad (25)$$

The function  $F_f$  obtained by substitution of (25) into (24) decreases with the sample size  $m$  as a  $O\left(\frac{\ln(m)}{\sqrt{m}}\right)$ . This convergence rate is that of the literature (for the margin loss function considered). The dependence on the number  $C$  of categories is a  $O\left(C^{\frac{d_{\mathcal{G}, C}}{2}} \ln(C)\right)$ , implying that it is always sublinear except in the worst case  $d_{\mathcal{G}, C} = 2$ . At last, the dependence on  $\frac{1}{\gamma}$  is a  $O\left(\left(\frac{1}{\gamma}\right)^{\frac{d_{\mathcal{G}, \gamma}}{2}}\right)$ . Even though the decrease of  $F_f$  with  $\gamma$  was expected, it calls for an explanation, since for a fixed value of the scale parameter  $\epsilon$ , the metric entropy and its upper bound (Formula (22)) increase with  $\gamma$ . The obvious reason is that  $F_i$  introduces a linear dependence of  $\epsilon$  on  $\gamma$  ( $\epsilon = \frac{\gamma}{2}$ ).

## 7.3 Guaranteed Risk for the Parameterized Truncated Hinge Loss

In the family of Lipschitz continuous margin loss functions, the option of choice is the parameterized truncated hinge loss  $\phi_{2, \gamma}$  given by:

$$\forall t \in \mathbb{R}, \quad \phi_{2, \gamma}(t) = \mathbf{1}_{\{t \leq 0\}} + \left(1 - \frac{t}{\gamma}\right) \mathbf{1}_{\{t \in (0, \gamma]\}}.$$

For this loss function, the best instance of Inequality (1) is provided by Theorem 5 in Guermeur (2017). It corresponds to  $L_* = L_\gamma$  and the analytical expression of function  $F_i$  is:

$$F_i(m, \gamma, \delta, \text{cap}(\rho_{\mathcal{G}, \gamma})) = \frac{2}{\gamma} R_m(\rho_{\mathcal{G}, \gamma}) + \sqrt{\frac{\ln(\frac{1}{\delta})}{2m}}. \quad (26)$$

Formula (26) calls for a discussion. On the one hand, a non trivial expression for the function  $F_f$  should decrease with  $\gamma$ . On the other hand, the Rademacher complexity is increasing with this parameter. This implies that we are looking for an upper bound on  $R_m(\rho_{\mathcal{G}, \gamma})$  increasing at most linearly with  $\gamma$ . In accordance with both graphs of transitions (Figures 1 and 2), this capacity measure is upper bounded as a function of the  $L_2$ -norm metric entropy by means of Dudley's chaining method. We use the following formula, whose degrees of freedom can be exploited to optimize the dependence on the basic parameters.

**Theorem 2** (Theorem 9 in Guermeur, 2017). *Let  $\mathcal{F}$  be a class of bounded real-valued functions on  $\mathcal{T}$ . For  $n \in \mathbb{N}^*$ , let  $\mathbf{t}_n \in \mathcal{T}^n$  and let  $\text{diam}(\mathcal{F})$  be the diameter of  $\mathcal{F}$  with respect to the pseudo-metric  $d_{2, \mathbf{t}_n}$ . Let  $h$  be a positive and decreasing function on  $\mathbb{N}$  such that  $h(0) \geq \text{diam}(\mathcal{F})$ . Then for  $N \in \mathbb{N}^*$ ,*

$$\hat{R}_n(\mathcal{F}) \leq h(N) + 2 \sum_{j=1}^N (h(j) + h(j-1)) \sqrt{\frac{\ln(\mathcal{N}^{\text{int}}(h(j), \mathcal{F}, d_{2, \mathbf{t}_n}))}{n}}. \quad (27)$$

A substitution of Inequality (23) into (27) gives:

$$R_m(\rho_{\mathcal{G}, \gamma}) \leq h(N) + \frac{8}{3} \sqrt{\frac{F_1(C)}{m}} \sum_{j \in \mathcal{J}} \frac{h(j) + h(j-1)}{h(j)^{\frac{d_{\mathcal{G}, \gamma}}{2}}} \log_2 \left( (C-1) \left( \frac{4\gamma}{h(j)} \right)^5 \right) \quad (28)$$

where

$$F_1(C) = 28^{d_{\mathcal{G}, \gamma}} K_{\rho_{\mathcal{G}}} C^{d_{\mathcal{G}, C}}, \quad (29)$$

with  $\mathcal{J} = \{j \in \llbracket 1; N \rrbracket : h(j) \leq \gamma\}$ . With the last formula at hand, the derivation of the confidence interval amounts to studying the phase transitions highlighted by Theorem 18 in Mendelson (2003).

**Theorem 3.** *Let  $\mathcal{G}$  be a function class satisfying Hypothesis 1. The following statements hold true for every value of  $\gamma$  in  $(0, 1]$ .*

*If  $d_{\mathcal{G}, \gamma} \in (0, 2)$ , then*

$$R_m(\rho_{\mathcal{G}, \gamma}) \leq 12 \left( 1 + 2^{\frac{2}{2-d_{\mathcal{G}, \gamma}}} \right) \sqrt{\frac{F_1(C)}{m}} F_2(C) \gamma^{1-\frac{d_{\mathcal{G}, \gamma}}{2}},$$

*where  $F_1(C)$  is given by Equation (29) and  $F_2(C) = \ln((C-1)4^5) + 10^{\frac{1+\ln(2)}{2-d_{\mathcal{G}, \gamma}}}$ .*

If  $d_{\mathcal{G},\gamma} = 2$  and  $m \geq 2$ , then

$$R_m(\rho_{\mathcal{G},\gamma}) \leq \gamma \frac{\log_2(m)}{\sqrt{m}} + 8 \sqrt{\frac{F_1(C)}{m}} \left[ \log_2 \left( \frac{\sqrt{m}}{\log_2(m)} \right) \right] \log_2 \left( (C-1) \left( 4 \frac{\sqrt{m}}{\log_2(m)} \right)^5 \right).$$

At last, if  $d_{\mathcal{G},\gamma} > 2$  and  $m \geq 2$ , then

$$R_m(\rho_{\mathcal{G},\gamma}) \leq \gamma \left( \frac{\log_2(m)}{m} \right)^{\frac{1}{d_{\mathcal{G},\gamma}}} \times \left[ 1 + \frac{16}{3} \left( 1 + 2^{\frac{2}{d_{\mathcal{G},\gamma}-2}} \right) \left( \frac{1}{\gamma} \right)^{\frac{d_{\mathcal{G},\gamma}}{2}} \sqrt{\frac{F_1(C)}{\log_2(m)}} \log_2 \left( (C-1) \left( 4 \left( \frac{m}{\log_2(m)} \right)^{\frac{1}{d_{\mathcal{G},\gamma}}} \right)^5 \right) \right].$$

The dependence of  $F_f$  on  $m$  is that of Mendelson's formulas, except for a multiplicative factor  $\sqrt{\ln(m)}$  in the case when  $d_{\mathcal{G},\gamma} \geq 2$  (complex classifiers). The dependences on  $C$  and  $\frac{1}{\gamma}$  are the same as with the indicator margin loss function.

#### 7.4 Comparison with the Use of the State-of-the-Art Structural Results

Compared to the use of the structural result involving Rademacher complexities (Lemma 1), the obvious advantage of our approach is to allow to exploit a possible coupling between the component functions of the classifier to produce a sublinear dependence of  $F_f$  on the number  $C$  of categories. The comparison with the decomposition involving covering numbers (Lemma 2) produces a mixed result. On the one hand, whatever the choice of the margin loss function, the dependence of  $F_f$  on the inverse of the margin parameter  $\gamma$  is improved. The gain is a factor  $\ln\left(\frac{1}{\gamma}\right)$  with the indicator margin loss function and a factor  $\sqrt{\ln\left(\frac{1}{\gamma}\right)}$  with the parameterized truncated hinge loss. The only prize to pay occurs for this latter margin loss function, when  $d_{\mathcal{G},\gamma} \geq 2$  (complex classifiers). Then, the dependence on the sample size  $m$  increases by a factor  $\sqrt{\ln(m)}$ .

## 8 Conclusions

We have established that the combinatorial results involving the fat-shattering dimension of the class of margin functions  $\rho_{\mathcal{G}}$  can always be improved by replacing this dimension with the margin Graph dimension of the same class (Lemmas 4, 6 and 7). Currently, the gain is limited by the lack of a structural result specific to the margin Graph dimension (a structural result that would be significantly better than Lemma 3). Fortunately, the use of another  $\gamma$ - $\Psi$ -dimension, the margin Natarajan dimension, makes it possible to exploit basic features of the classifier of interest (Corollary 1 and Lemma 11) to derive useful structural

results. The major consequence is an improved dependence of the confidence interval of the guaranteed risk on the margin parameter  $\gamma$ . This holds true both with the  $L_\infty$ -norm (Inequality (25)) and the  $L_2$ -norm (Theorem 3). Except in the worst case  $d_{\mathcal{G},C} = 2$ , the dependence on the number  $C$  of categories is sublinear. The only drawback is that the convergence rate of the guaranteed risk associated with the parameterized truncated hinge loss can be worsened by a factor  $\sqrt{\ln(m)}$  when the class  $\mathcal{G}_0$  is complex (large value of  $d_{\mathcal{G},\gamma}$ ). The phenomenon is a direct consequence of the appearance of the logarithmic function of  $\frac{\gamma}{\epsilon}$  in the exponent of the right-hand side of Inequality (14) (compared to Formula 2 in Mendelson and Vershynin, 2003). Whether this term can be replaced with a logarithmic function of  $C$  only is an open question which is the subject of an ongoing research.

**Acknowledgements** The author would like to thank R. Vershynin for his explanations on the proof of Theorem 1 in Mendelson and Vershynin (2003). Thanks are also due to F. Lauer and T. Masini for carefully reading this manuscript. This work was partly funded by a CNRS research grant (PEPS).

## References

- N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler. Scale-sensitive dimensions, uniform convergence, and learnability. *Journal of the ACM*, 44(4):615–631, 1997.
- M. Anthony and P.L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, Cambridge, 1999.
- P.L. Bartlett. The sample complexity of pattern classification with neural networks: The size of the weights is more important than the size of the network. *IEEE Transactions on Information Theory*, 44(2):525–536, 1998.
- P.L. Bartlett and P.M. Long. Prediction, learning, uniform convergence, and scale-sensitive dimensions. *Journal of Computer and System Sciences*, 56(2):174–190, 1998.
- P.L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- P.L. Bartlett and J. Shawe-Taylor. Generalization performance of support vector machines and other pattern classifiers. In B. Schölkopf, C.J.C. Burges, and A. Smola, editors,

- Advances in Kernel Methods - Support Vector Learning*, chapter 4, pages 43–54. The MIT Press, Cambridge, MA, 1999.
- S. Ben-David, N. Cesa-Bianchi, D. Haussler, and P.M. Long. Characterizations of learnability for classes of  $\{0, \dots, n\}$ -valued functions. *Journal of Computer and System Sciences*, 50(1):74–86, 1995.
- A. Berlinet and C. Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Kluwer Academic Publishers, Boston, 2004.
- U. Doğan, T. Glasmachers, and C. Igel. A unified view on multi-class support vector classification. *Journal of Machine Learning Research*, 17(45):1–32, 2016.
- H.H. Duan. Bounding the fat shattering dimension of a composition function class built using a continuous logic connective. *The Waterloo Mathematics Review*, 2(1):1–21, 2012.
- R.M. Dudley. The sizes of compact subsets of Hilbert space and continuity of Gaussian processes. *Journal of Functional Analysis*, 1(3):290–330, 1967.
- R.M. Dudley, E. Giné, and J. Zinn. Uniform and universal Glivenko-Cantelli classes. *Journal of Theoretical Probability*, 4(3):485–510, 1991.
- Y. Guermeur. VC theory of large margin multi-category classifiers. *Journal of Machine Learning Research*, 8:2551–2594, 2007.
- Y. Guermeur.  $L_p$ -norm Sauer-Shelah lemma for margin multi-category classifiers. *Journal of Computer and System Sciences*, 89:450–473, 2017.
- M.J. Kearns and R.E. Schapire. Efficient distribution-free learning of probabilistic concepts. *Journal of Computer and System Sciences*, 48(3):464–497, 1994.
- A.N. Kolmogorov and V.M. Tihomirov.  $\epsilon$ -entropy and  $\epsilon$ -capacity of sets in functional spaces. *American Mathematical Society Translations, series 2*, 17:277–364, 1961.
- A. Kontorovich and R. Weiss. Maximum margin multiclass nearest neighbors. In *ICML’14*, 2014.
- Y. Lei, U. Doğan, A. Binder, and M. Kloft. Multi-class SVMs: From tighter data-dependent generalization bounds to novel algorithms. In *NIPS 28*, pages 2026–2034, 2015.
- A. Maurer. A vector-contraction inequality for Rademacher complexities. In *ALT’16*, pages 3–17, 2016.

- S. Mendelson. A few notes on statistical learning theory. In S. Mendelson and A.J. Smola, editors, *Advanced Lectures on Machine Learning*, chapter 1, pages 1–40. Springer-Verlag, Berlin, Heidelberg, New York, 2003.
- S. Mendelson and R. Vershynin. Entropy and the combinatorial dimension. *Inventiones mathematicae*, 152:37–55, 2003.
- M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of Machine Learning*. The MIT Press, Cambridge, MA, second edition, 2018.
- K. Musayeva, F. Lauer, and Y. Guermeur. Rademacher complexity and generalization performance of multi-category margin classifiers. *Neurocomputing*, 342:6–15, 2019.
- D. Pollard. *Convergence of Stochastic Processes*. Springer-Verlag, New York, 1984.
- M. Rudelson and R. Vershynin. Combinatorics of random processes and sections of convex bodies. *Annals of Mathematics*, 164(2):603–648, 2006.
- M. Talagrand. Vapnik-Chervonenkis type conditions and uniform Donsker classes of functions. *The Annals of Probability*, 31(3):1565–1582, 2003.
- M. Talagrand. *Upper and Lower Bounds for Stochastic Processes: Modern Methods and Classical Problems*. Springer-Verlag, Berlin Heidelberg, 2014.
- A.W. van der Vaart and J.A. Wellner. *Weak Convergence and Empirical Processes, With Applications to Statistics*. Springer Series in Statistics. Springer-Verlag, New York, 1996.
- V.N. Vapnik and A.Ya. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, XVI(2):264–280, 1971.
- G. Wahba. Multivariate function and operator estimation, based on smoothing splines and reproducing kernels. In M. Casdagli and S. Eubank, editors, *Nonlinear Modeling and Forecasting, SFI Studies in the Sciences of Complexity*, volume XII, pages 95–112. Addison-Wesley, 1992.



## A Proofs of the Connections Between the Scale-Sensitive Combinatorial Dimensions

The proofs in this appendix and the following ones make use of a standard convention: a function class  $\mathcal{F}$  is said to  $\gamma$ -N-shatter a triplet  $(s_{\mathcal{Z}^n}, \mathbf{b}_n, \mathbf{c}_n)$  if  $\mathcal{F}$   $\gamma$ -N-shatters  $s_{\mathcal{Z}^n}$  and  $(\mathbf{b}_n, \mathbf{c}_n)$  is a witness to this shattering. The corresponding convention for the margin Graph dimension is also used. The proof of Lemma 4 is the following one.

*Proof.* Let  $f$  be any function in  $\mathcal{F}$  and  $z = (x, y) \in \mathcal{Z}$ . Let  $\gamma \in (0, 1]$ ,  $\epsilon \in (0, \frac{\gamma}{2}]$  and  $b \in [\epsilon, \gamma - \epsilon]$ . Then,

$$f_\gamma(z) - b \geq \epsilon \implies f_\gamma(z) \geq 2\epsilon > 0 \implies f(z) \geq f_\gamma(z).$$

Consequently,

$$f_\gamma(z) - b \geq \epsilon \implies f(z) - b \geq \epsilon. \quad (30)$$

Suppose that  $f_\gamma(z) = 0$  (implying that  $f(z) \leq 0$ ). Then it results from Equation (2) that

$$-f_\gamma(z) \leq \max_{k \neq y} f(x, k) \leq -f(z).$$

Suppose now that  $f_\gamma(z) > 0$ . Then,

$$-f_\gamma(z) + b \geq \epsilon \implies f_\gamma(z) \leq \gamma - 2\epsilon < \gamma.$$

Thus,  $f_\gamma(z)$  belongs to the open interval  $(0, \gamma)$ , so that  $f(z) = f_\gamma(z)$ . Consequently,

$$-f_\gamma(z) = \max_{k \neq y} f(x, k) = -f(z).$$

To sum up, for all possible values of  $f_\gamma(z)$  (positive or null),

$$-f_\gamma(z) + b \geq \epsilon \implies \max_{k \neq y} f(x, k) + b \geq \epsilon \implies -f(z) + b \geq \epsilon. \quad (31)$$

Formula (9) directly springs from Formulas (30) and (31).  $\square$

The proof of Lemma 5 is the following one.

*Proof.* The left-hand side inequality of Formula (10) is obvious. Let us set  $d_G = \gamma$ -G-dim  $(\mathcal{F})$  and  $d_N = \gamma$ -N-dim  $(\mathcal{F})$ . The right-hand side inequality of Formula (10) is trivially true for  $d_G = 0$ . Thus, we prove it under the assumption that  $d_G \geq 1$ . Let  $\tilde{\mathcal{F}}$  be any subset of  $\mathcal{F}$  of cardinality  $2^{d_G}$  that  $\gamma$ -G-shatters a subset  $s_{\mathcal{Z}^{d_G}}$  of  $\mathcal{Z}$  of cardinality  $d_G$ . The proof rests on the derivation of a separating tree of  $\tilde{\mathcal{F}}$ . For notational simplicity, we set

$s_{\mathcal{Z}^{d_G}} = \{z_i : 1 \leq i \leq d_G\}$ . A subset of  $s_{\mathcal{Z}^{d_G}}$  of cardinality  $n \in \llbracket 1; d_N \rrbracket$  is denoted by  $s'_{\mathcal{Z}^n} = \{z'_i : 1 \leq i \leq n\}$ , with the convention

$$\forall (i, j) : 1 \leq i < j \leq n, (z'_i, z'_j) = (z_v, z_w) \implies 1 \leq v < w \leq d_G.$$

Let the vector  $\mathbf{b}_{d_G} = (b_i)_{1 \leq i \leq d_G} \in \mathbb{R}_+^{d_G}$  be a witness to the  $\gamma$ -G-shattering of  $s_{\mathcal{Z}^{d_G}}$  by  $\tilde{\mathcal{F}}$  and let  $(s'_{\mathcal{Z}^n}, \mathbf{b}'_n, \mathbf{c}'_n)$  be any triplet such that

$$\begin{cases} s'_{\mathcal{Z}^n} \subset s_{\mathcal{Z}^{d_G}} \\ \mathbf{b}'_n \in \mathbb{R}_+^n : \forall i \in \llbracket 1; n \rrbracket, z'_i = z_j \implies b'_i = b_j \\ \mathbf{c}'_n \in \mathcal{Y}^n : \forall i \in \llbracket 1; n \rrbracket, c'_i \in \mathcal{Y} \setminus \{y'_i\} \end{cases} .$$

For every subset  $\bar{\mathcal{F}}$  of  $\tilde{\mathcal{F}}$ , let  $s(\bar{\mathcal{F}})$  denote the number of triplets  $(s'_{\mathcal{Z}^n}, \mathbf{b}'_n, \mathbf{c}'_n)$   $\gamma$ -N-shattered by  $\bar{\mathcal{F}}$ . Combinatorics produces

$$s(\tilde{\mathcal{F}}) \leq \sum_{n=1}^{d_N} \binom{d_G}{n} (C-1)^n,$$

which gives birth to a handy formula thanks to a well-known computation (see for instance the proof of Corollary 3.18 in Mohri et al., 2018):

$$s(\tilde{\mathcal{F}}) \leq \left( \frac{(C-1)ed_G}{d_N} \right)^{d_N} - 1. \quad (32)$$

The derivation of the separating tree of  $\tilde{\mathcal{F}}$  provides us with a lower bound on  $s(\tilde{\mathcal{F}})$ . Let  $\bar{\mathcal{F}}$  be any of its nodes such that  $|\bar{\mathcal{F}}| \geq 2$  (inner node). Its two sons,  $\bar{\mathcal{F}}_+$  and  $\bar{\mathcal{F}}_-$ , are built as follows. Split  $\bar{\mathcal{F}}$  arbitrarily into  $\left\lfloor \frac{|\bar{\mathcal{F}}|}{2} \right\rfloor$  pairs (with possibly a function remaining alone). For each pair  $(f, f')$ , find  $z_i \in s_{\mathcal{Z}^{d_G}}$  such that

$$\begin{cases} f(z_i) - b_i \geq \gamma \\ \max_{k \neq y_i} f'(x_i, k) + b_i \geq \gamma \end{cases}$$

or vice versa. By the pigeonhole principle, the same example is picked for at least  $\left\lceil \left\lfloor \frac{|\bar{\mathcal{F}}|}{2} \right\rfloor \frac{1}{d_G} \right\rceil$  pairs. Let  $z_{i_0}$  be such an example, and let  $(f_+, f_-)$  denote the corresponding pairs, whose components are reordered (when needed) so that

$$\begin{cases} f_+(z_{i_0}) - b_{i_0} \geq \gamma \\ \max_{k \neq y_{i_0}} f_-(x_{i_0}, k) + b_{i_0} \geq \gamma \end{cases} .$$

$\bar{\mathcal{F}}_+$  is the set of the functions  $f_+$ . Once more by the pigeonhole principle, there exists a value  $c_{i_0} \in \mathcal{Y} \setminus \{y_{i_0}\}$  such that at least  $\left\lceil \left\lceil \left\lceil \frac{|\bar{\mathcal{F}}|}{2} \right\rceil \frac{1}{d_G} \right\rceil \frac{1}{C-1} \right\rceil$  functions  $f_-$  satisfy

$$f_-(x_{i_0}, c_{i_0}) + b_{i_0} \geq \gamma.$$

Let  $\bar{\mathcal{F}}_-$  be the set of these functions  $f_-$ . The two sons  $\bar{\mathcal{F}}_+$  and  $\bar{\mathcal{F}}_-$  of  $\bar{\mathcal{F}}$  have been built in such a way that  $|\bar{\mathcal{F}}_+| \geq \frac{|\bar{\mathcal{F}}|}{3d_G}$ ,  $|\bar{\mathcal{F}}_-| \geq \frac{|\bar{\mathcal{F}}|}{3d_G(C-1)}$  and

$$\begin{cases} \forall f_+ \in \bar{\mathcal{F}}_+, & f_+(z_{i_0}) - b_{i_0} \geq \gamma \\ \forall f_- \in \bar{\mathcal{F}}_-, & f_-(x_{i_0}, c_{i_0}) + b_{i_0} \geq \gamma \end{cases}. \quad (33)$$

Since  $\bar{\mathcal{F}}_+ \cup \bar{\mathcal{F}}_- \subset \bar{\mathcal{F}}$ , any triplet  $\gamma$ -N-shattered by either  $\bar{\mathcal{F}}_+$  or  $\bar{\mathcal{F}}_-$  is also  $\gamma$ -N-shattered by  $\bar{\mathcal{F}}$ . Furthermore, according to (33), if the triplet  $(s'_{\mathcal{Z}^n}, \mathbf{b}'_n, \mathbf{c}'_n)$  is  $\gamma$ -N-shattered by both  $\bar{\mathcal{F}}_+$  and  $\bar{\mathcal{F}}_-$ , then  $\bar{\mathcal{F}}$  also  $\gamma$ -N-shatters the triplet  $(s''_{\mathcal{Z}^{n+1}}, \mathbf{b}''_{n+1}, \mathbf{c}''_{n+1})$  such that  $s''_{\mathcal{Z}^{n+1}} = s'_{\mathcal{Z}^n} \cup \{z_{i_0}\}$ , the vector  $\mathbf{b}''_{n+1}$  is deduced from  $\mathbf{b}'_n$  by inserting the component  $b_{i_0}$ , and the vector  $\mathbf{c}''_{n+1}$  is deduced from  $\mathbf{c}'_n$  by inserting the component  $c_{i_0}$ . Clearly, neither  $\bar{\mathcal{F}}_+$  nor  $\bar{\mathcal{F}}_-$   $\gamma$ -N-shatters the triplet  $(s''_{\mathcal{Z}^{n+1}}, \mathbf{b}''_{n+1}, \mathbf{c}''_{n+1})$ , simply because (contrary to  $\bar{\mathcal{F}}$ ) they do not  $\gamma$ -N-shatter the triplet  $(\{z_{i_0}\}, b_{i_0}, c_{i_0})$ . A synthesis of the different cases produces:

$$s(\bar{\mathcal{F}}) \geq s(\bar{\mathcal{F}}_+) + s(\bar{\mathcal{F}}_-) + 1. \quad (34)$$

To further bound from below  $s(\bar{\mathcal{F}})$ , we introduce the function  $\ell$  which returns the number of leaves of the (sub)tree whose root is its argument. Using (34), the following simple connection between the two functions can be proved by induction:

$$s(\bar{\mathcal{F}}) \geq \ell(\bar{\mathcal{F}}) - 1. \quad (35)$$

Once more by induction, we now establish that any node  $\bar{\mathcal{F}}$  (even a leaf) satisfies:

$$\ell(\bar{\mathcal{F}}) \geq |\bar{\mathcal{F}}|^{\frac{1}{\log_2(3\sqrt{C-1}d_G)}}. \quad (36)$$

Inequality (36) is obviously true for the leaves (for which  $\ell(\bar{\mathcal{F}}) = |\bar{\mathcal{F}}| = 1$ ). Suppose now that it holds true for the two sons  $\bar{\mathcal{F}}_+$  and  $\bar{\mathcal{F}}_-$  of the inner node  $\bar{\mathcal{F}}$ . Then,

$$\begin{aligned}
\ell(\tilde{\mathcal{F}}) &= \ell(\tilde{\mathcal{F}}_+) + \ell(\tilde{\mathcal{F}}_-) \\
&\geq \left( \frac{|\tilde{\mathcal{F}}|}{3d_G} \right)^{\frac{1}{\log_2(3\sqrt{C-1}d_G)}} + \left( \frac{|\tilde{\mathcal{F}}|}{3d_G(C-1)} \right)^{\frac{1}{\log_2(3\sqrt{C-1}d_G)}} \\
&= \frac{1}{2} \left( \left( \sqrt{C-1} \right)^{\frac{1}{\log_2(3\sqrt{C-1}d_G)}} + \left( \sqrt{C-1} \right)^{-\frac{1}{\log_2(3\sqrt{C-1}d_G)}} \right) |\tilde{\mathcal{F}}|^{\frac{1}{\log_2(3\sqrt{C-1}d_G)}} \\
&\geq \frac{1}{2} \min_{t \in \mathbb{R}_+^*} \left( t + \frac{1}{t} \right) |\tilde{\mathcal{F}}|^{\frac{1}{\log_2(3\sqrt{C-1}d_G)}} \\
&= |\tilde{\mathcal{F}}|^{\frac{1}{\log_2(3\sqrt{C-1}d_G)}}.
\end{aligned}$$

Since  $|\tilde{\mathcal{F}}| = 2^{d_G}$ , we thus get for the whole tree:

$$\ell(\tilde{\mathcal{F}}) \geq 2^{\frac{d_G}{\log_2(3\sqrt{C-1}d_G)}}. \quad (37)$$

A substitution of (37) into (35) provides the lower bound on  $s(\tilde{\mathcal{F}})$  announced:

$$s(\tilde{\mathcal{F}}) \geq 2^{\frac{d_G}{\log_2(3\sqrt{C-1}d_G)}} - 1. \quad (38)$$

Combining (38) and the upper bound, (32), gives by transitivity:

$$\begin{aligned}
d_G &\leq d_N \log_2 \left( \frac{(C-1)ed_G}{d_N} \right) \log_2 \left( 3\sqrt{C-1}d_G \right) \\
&\leq \frac{1}{\ln^2(2)} d_N \ln \left( F(C) \frac{d_G}{d_N} \right) \ln(F(C)d_G),
\end{aligned} \quad (39)$$

where  $F(C) = e(C-1)$ . To bound from above the right-hand side of Inequality (39), we resort to the following statement:

$$\forall (u, u_0) \in [1, +\infty)^2, \ln(u) \leq 2u_0 u^{\frac{1}{4u_0}}, \quad (40)$$

with  $u_0 = \ln(F(C))$ . We then obtain

$$\begin{cases} \ln \left( F(C) \frac{d_G}{d_N} \right) \leq 2e^{\frac{1}{4}} \ln(F(C)) \left( \frac{d_G}{d_N} \right)^{\frac{1}{4\ln(F(C))}} \\ \ln(F(C)d_G) \leq 2e^{\frac{1}{4}} \ln(F(C)) d_G^{\frac{1}{4\ln(F(C))}} \end{cases}.$$

By substitution into (39),

$$\begin{aligned}
d_G &\leq \frac{4\sqrt{e}}{\ln^2(2)} \ln^2(F(C)) d_G^{\frac{1}{2\ln(F(C))}} d_N^{\frac{4\ln(F(C))-1}{4\ln(F(C))}} \\
&\leq \left(\frac{4\sqrt{e}}{\ln^2(2)}\right)^{\frac{2\ln(F(C))}{2\ln(F(C))-1}} (\ln(F(C)))^{\frac{4\ln(F(C))}{2\ln(F(C))-1}} d_N^{\frac{4\ln(F(C))-1}{4\ln(F(C))-2}} \\
&= \left(\frac{4\sqrt{e}}{\ln^2(2)}\right)^{\frac{2\ln(F(C))}{2\ln(F(C))-1}} (\ln(F(C)))^{\frac{2}{2\ln(F(C))-1}} \ln^2(2) \log_2^2(F(C)) d_N^{\frac{4\ln(F(C))-1}{4\ln(F(C))-2}} \\
&< 32 \log_2^2(F(C)) d_N^{1+\frac{1}{4\ln(F(C))-2}}.
\end{aligned}$$

□

## B Proofs of the Combinatorial Results

This appendix gathers the proofs of the four new combinatorial results. It starts with three lemmas which are common to all proofs.

### B.1 Shared Technical Lemmas

Each of the combinatorial results in the literature is built upon a basic lemma that involves two (possibly identical) function classes whose codomains are discrete. The domain and codomain of the first one are finite sets, so that its cardinality is also finite. This cardinality is upper bounded in terms of a combinatorial dimension of the second function class. In the case of margin classifiers, the combinatorial dimension of the basic lemma is a variant of the scale-sensitive dimension of the combinatorial result, variant designed to take benefit from the aforementioned restrictions. The first capacity measure of this kind is a variant of the  $\gamma$ -dimension: the strong dimension (Definition 3.1 in Alon et al., 1997). The strong  $\Psi$ -dimensions extend the  $\gamma$ - $\Psi$ -dimensions according to the same principle.

**Definition 12** (Strong  $\Psi$ -dimensions). *Let  $\mathcal{F}$  be a function class defined as in Definition 8. Suppose further that the functions in  $\mathcal{F}$  take their values in  $\mathbb{Z}$ . Let  $\Psi$  be a family of mappings from  $\mathcal{Y}$  into  $\{-1, 0, 1\}$ . A subset  $s_{\mathcal{Z}^n} = \{z_i = (x_i, y_i) : 1 \leq i \leq n\}$  of  $\mathcal{Z}$  is said to be strongly  $\Psi$ -shattered by  $\mathcal{F}$  if there is a vector  $\boldsymbol{\psi}_n = (\psi^{(i)})_{1 \leq i \leq n} \in \Psi^n$  satisfying  $(\psi^{(i)}(y_i))_{1 \leq i \leq n} = \mathbf{1}_n$ , and a vector  $\mathbf{b}_n = (b_i)_{1 \leq i \leq n} \in \mathbb{N}^n$  such that, for every vector  $\mathbf{s}_n = (s_i)_{1 \leq i \leq n} \in \{-1, 1\}^n$ , there is a function  $f_{\mathbf{s}_n} \in \mathcal{F}$  satisfying*

$$\forall i \in \llbracket 1; n \rrbracket, \quad s_i \left( s_i \max_{\{k: \psi^{(i)}(k)=s_i\}} f_{\mathbf{s}_n}(x_i, k) - b_i \right) \geq 1.$$

The strong  $\Psi$ -dimension of  $\mathcal{F}$ , denoted by  $S\text{-}\Psi\text{-dim}(\mathcal{F})$ , is the maximal cardinality of a subset of  $\mathcal{Z}$  strongly  $\Psi$ -shattered by  $\mathcal{F}$ , if such maximum exists. Otherwise,  $\mathcal{F}$  is said to have infinite strong  $\Psi$ -dimension.

In what follows, the finiteness of the domain is simply obtained by application of a restriction to (projection on) an appropriately chosen set of data points. The discretization of the codomain results from the application of the following operator.

**Definition 13** ( $\eta$ -discretization operator, Definition 33 in Guermeur, 2007). Let  $\mathcal{F} \subset \mathbb{R}^{\mathcal{T}}$ . For  $\eta \in \mathbb{R}_+^*$ , define the  $\eta$ -discretization as an operator on  $\mathcal{F}$  such that:

$$\begin{aligned} (\cdot)^{(\eta)} : \mathcal{F} &\longrightarrow \mathcal{F}^{(\eta)} \\ f &\mapsto f^{(\eta)} \end{aligned}$$

$$\forall t \in \mathcal{T}, \quad f^{(\eta)}(t) = \text{sign}(f(t)) \cdot \left\lfloor \frac{|f(t)|}{\eta} \right\rfloor.$$

The transitions from continuous functions to discrete ones and back are obtained by application of the two following lemmas.

**Lemma 12.** Let  $\mathcal{F}$  be a class of functions on  $\mathcal{T}$  taking nonnegative values ( $\mathcal{F} \subset \mathbb{R}_+^{\mathcal{T}}$ ). For  $n \in \mathbb{N}^*$ , let  $\mathbf{t}_n = (t_i)_{1 \leq i \leq n} \in \mathcal{T}^n$ . Let  $N$  be a positive integer. For every  $\epsilon \in \mathbb{R}_+^*$  and every  $\eta \in \left(0, \frac{\epsilon}{N+1}\right]$ ,

$$\forall (f, f') \in \mathcal{F}^2, \quad d_{2, \mathbf{t}_n}(f, f') \geq \epsilon \implies d_{2, \mathbf{t}_n}(f^{(\eta)}, f'^{(\eta)}) \geq N, \quad (41)$$

with the consequence that if the subset  $\bar{\mathcal{F}}$  of  $\mathcal{F}$  is  $\epsilon$ -separated with respect to the pseudo-metric  $d_{2, \mathbf{t}_n}$ , then it is in bijection with the subset  $\bar{\mathcal{F}}^{(\eta)}$  of  $\mathcal{F}^{(\eta)}$ , which is  $N$ -separated with respect to the same pseudo-metric. Similarly, for every  $\epsilon \in \mathbb{R}_+^*$  and every  $\eta \in \left(0, \frac{\epsilon}{2}\right]$ ,

$$\mathcal{M}(\epsilon, \mathcal{F}, d_{\infty, \mathbf{t}_n}) \leq \mathcal{M}\left(2, \mathcal{F}^{(\eta)}, d_{\infty, \mathbf{t}_n}\right). \quad (42)$$

*Proof.* For  $f \in \mathcal{F}$  and  $i \in \llbracket 1; n \rrbracket$ , let us denote the Euclidean division of  $f(t_i)$  by  $\eta$  as follows:

$$\forall i \in \llbracket 1; n \rrbracket, \quad f(t_i) = \eta f^{(\eta)}(t_i) + r_i.$$

With the notation introduced above,

$$d_{2, \mathbf{t}_n}(f, f')^2 = \frac{1}{n} \sum_{i=1}^n \left( \eta \left( f^{(\eta)}(t_i) - f'^{(\eta)}(t_i) \right) + r_i - r'_i \right)^2.$$

For  $i \in \llbracket 1; n \rrbracket$ , let  $\delta_i = \left| f^{(\eta)}(t_i) - f'^{(\eta)}(t_i) \right|$ .

$$\begin{aligned} (d_{2, \mathbf{t}_n}(f, f') \geq \epsilon) \text{ and } \left( \eta \in \left( 0, \frac{\epsilon}{N+1} \right] \right) &\implies \left( \frac{1}{n} \sum_{i=1}^n (\eta \delta_i + |r_i - r'_i|)^2 \right)^{\frac{1}{2}} \geq \epsilon \\ &\implies \left( \frac{1}{n} \sum_{i=1}^n (\delta_i + 1)^2 \right)^{\frac{1}{2}} \geq \frac{\epsilon}{\eta} \\ &\implies \left( \frac{1}{n} \sum_{i=1}^n (\delta_i + 1)^2 \right)^{\frac{1}{2}} \geq N+1 \quad (43) \end{aligned}$$

$$\begin{aligned} &\implies \left( \frac{1}{n} \sum_{i=1}^n \delta_i^2 \right)^{\frac{1}{2}} + 1 \geq N+1 \quad (44) \\ &\implies d_{2, \mathbf{t}_n}(f^{(\eta)}, f'^{(\eta)}) \geq N, \end{aligned}$$

where the transition from (43) to (44) is provided by the triangle inequality. To sum up, we have established (41), i.e., the part of the lemma dealing with the  $L_2$ -norm. To prove (42), it is enough to observe that

$$\begin{aligned} f(t) - f'(t) \geq \epsilon &\iff \frac{\epsilon}{2} \left( f^{(\frac{\epsilon}{2})}(t) - f'^{(\frac{\epsilon}{2})}(t) \right) + r - r' \geq \epsilon \\ &\implies \frac{\epsilon}{2} \left( f^{(\frac{\epsilon}{2})}(t) - f'^{(\frac{\epsilon}{2})}(t) + 1 \right) > \epsilon \\ &\implies f^{(\frac{\epsilon}{2})}(t) - f'^{(\frac{\epsilon}{2})}(t) \geq 2. \end{aligned}$$

□

**Lemma 13.** *Let  $\mathcal{F}$  be a function class defined as in Definition 8. For every  $\eta \in \mathbb{R}_+^*$  and every  $\epsilon \in (0, \frac{\eta}{2}]$ ,*

$$S\text{-}\Psi\text{-dim}(\mathcal{F}^{(\eta)}) \leq \epsilon\text{-}\Psi\text{-dim}(\mathcal{F}). \quad (45)$$

*Proof.* To prove (45), it is enough to notice that any set  $s_{\mathcal{Z}^n}$  strongly  $\Psi$ -shattered by  $\mathcal{F}^{(\eta)}$  according to the vector  $\mathbf{b}_n = (b_i)_{1 \leq i \leq n} \in \mathbb{N}^n$  is also  $\frac{\eta}{2}$ - $\Psi$ -shattered by  $\mathcal{F}$  according to  $\mathbf{b}'_n = (\eta(b_i + \frac{1}{2}))_{1 \leq i \leq n} \in \mathbb{R}_+^n$ . □

**Lemma 14.** *Let  $\mathcal{F}$  be a function class defined as in Definition 8. Suppose that there exist  $(f, f') \in \mathcal{F}^2$ ,  $\gamma \in (0, 1]$ ,  $\eta \in (0, \frac{\gamma}{2}]$ , and  $z = (x, y) \in \mathcal{Z}$  such that*

$$f_\gamma^{(\eta)}(z) - f'_\gamma^{(\eta)}(z) \geq 2,$$

where  $f_\gamma^{(\eta)} = (\pi_\gamma \circ f)^{(\eta)}$  and  $f'_\gamma^{(\eta)}$  is defined accordingly. For every  $b \in \llbracket f'_\gamma^{(\eta)}(z) + 1; f_\gamma^{(\eta)}(z) - 1 \rrbracket$  and  $c \in \operatorname{argmax}_{k \neq y} f'^{(\eta)}(x, k)$ ,

1. the set  $\{f_\gamma^{(\eta)}, f_\gamma^{\prime(\eta)}\}$  strongly shatters the pair  $(\{z\}, b)$ ;
2. the set  $\{f^{(\eta)}, f^{\prime(\eta)}\}$  strongly  $G$ -shatters the same pair;
3. the set  $\{f^{(\eta)}, f^{\prime(\eta)}\}$  strongly  $N$ -shatters the triplets  $(\{z\}, b, c)$ .

*Proof.* 1. The first assertion directly springs from the definition of the fat-shattering dimension.

2. The second assertion can be derived from the first one following the line of reasoning of the proof of Lemma 4.

3. The third assertion directly springs from the second one and the definition of the margin Natarajan dimension. □

## B.2 Margin Graph Dimension - Uniform Convergence Norm

The proof of Lemma 6 borrows from the proofs of classical results, including the two state-of-the-art combinatorial results: Lemma 3.5 in Alon et al. (1997) and Theorem 1 in Mendelson and Vershynin (2003). Central in this proof is the following basic combinatorial result.

**Lemma 15.** *Let  $\mathcal{F}$  be a function class defined as in Definition 8. For every  $s_{\mathcal{Z}^n} = \{z_i = (x_i, y_i) : 1 \leq i \leq n\} \subset \mathcal{Z}$ ,  $\gamma \in (0, 1]$  and  $\eta \in (0, \frac{\gamma}{2}]$ , if  $\mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}}$  is 2-separated in the metric  $d_{\infty, \mathcal{Z}^n}$ , then*

$$\left| \mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}} \right| \leq (3M_\gamma n)^{\log_2(\Sigma)}, \quad (46)$$

where  $\Sigma = \sum_{u=0}^{d_G} \binom{n}{u} M_\gamma^u$  with  $M_\gamma = \lfloor \frac{\gamma}{\eta} \rfloor$  and  $d_G = S\text{-}G\text{-dim}(\mathcal{F}^{(\eta)})$ .

*Proof.* Notice first that Inequality (46) is trivially true for  $\left| \mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}} \right| = 1$ . Indeed, the minimal value of its right-hand side, corresponding to  $d_G = 0$ , is 1. Thus, the rest of the proof makes use of the restriction  $\left| \mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}} \right| \geq 2$ . A direct consequence is that according to Lemma 14,  $d_G \geq 1$ . A subset of  $s_{\mathcal{Z}^n}$  of cardinality  $u \in \llbracket 1; n \rrbracket$  is denoted by  $s'_{\mathcal{Z}^u} = \{z'_i : 1 \leq i \leq u\}$ , with the convention

$$\forall (i, j) : 1 \leq i < j \leq u, (z'_i, z'_j) = (z_v, z_w) \implies 1 \leq v < w \leq n.$$

For every subset  $\bar{\mathcal{F}}$  of  $\mathcal{F}$ , let  $s(\bar{\mathcal{F}}^{(\eta)})$  denote the number of pairs  $(s'_{\mathcal{Z}^u}, \mathbf{b}'_u)$  with  $s'_{\mathcal{Z}^u} \subset s_{\mathcal{Z}^n}$  and  $\mathbf{b}'_u \in \llbracket 1; M_\gamma - 1 \rrbracket^u$  strongly  $G$ -shattered by  $\bar{\mathcal{F}}^{(\eta)}$  (the convention above has been



introduced to avoid handling duplicates). Let  $\tilde{\mathcal{F}}$  be any subset of  $\mathcal{F}$  such that  $\tilde{\mathcal{F}}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}} = \mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}}$  and  $|\tilde{\mathcal{F}}| = \left| \mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}} \right|$  ( $\tilde{\mathcal{F}}$  and  $\mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}}$  are in bijection). Since  $\tilde{\mathcal{F}}^{(\eta)} \subset \mathcal{F}^{(\eta)}$  and  $d_G \geq 1$ , combinatorics gives:

$$\begin{aligned} s\left(\tilde{\mathcal{F}}^{(\eta)}\right) &\leq s\left(\mathcal{F}^{(\eta)}\right) \\ &\leq \sum_{u=1}^{d_G} \binom{n}{u} M_\gamma^u = \Sigma - 1. \end{aligned} \quad (47)$$

In order to derive a lower bound on  $s\left(\tilde{\mathcal{F}}^{(\eta)}\right)$ , we build a 2-separating tree of  $\tilde{\mathcal{F}}$ , i.e., a 2-separating tree of  $\mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}}$  (see Definition 3.4 in Rudelson and Vershynin, 2006). Let  $\bar{\mathcal{F}}$  be one of its nodes such that  $|\bar{\mathcal{F}}| \geq 2$  (inner node). Its two sons,  $\bar{\mathcal{F}}_+$  and  $\bar{\mathcal{F}}_-$ , are built as follows. Split  $\bar{\mathcal{F}}$  arbitrarily into  $\left\lfloor \frac{|\bar{\mathcal{F}}|}{2} \right\rfloor$  pairs (with possibly a function remaining alone). For each pair  $(f, f')$ , find  $z_i \in s_{\mathcal{Z}^n}$  such that  $\left| f_\gamma^{(\eta)}(z_i) - f'_{\gamma}{}^{(\eta)}(z_i) \right| \geq 2$ . By the pigeonhole principle, the same example is picked for at least  $\left\lfloor \left\lfloor \frac{|\bar{\mathcal{F}}|}{2} \right\rfloor \frac{1}{n} \right\rfloor$  pairs. Let  $z_{i_0}$  be such an example, and let  $(f_+, f_-)$  denote the corresponding pairs, whose components are reordered (when needed) so that

$$f_{+, \gamma}^{(\eta)}(z_{i_0}) > f_{-, \gamma}^{(\eta)}(z_{i_0}).$$

Among the functions  $f_{+, \gamma}^{(\eta)}$ , at least  $\left\lfloor \left\lfloor \left\lfloor \frac{|\bar{\mathcal{F}}|}{2} \right\rfloor \frac{1}{n} \right\rfloor \frac{1}{M_\gamma - 1} \right\rfloor$  take the same value at  $z_{i_0}$ . Let  $v(z_{i_0})$  be such a value. We define  $\bar{\mathcal{F}}_+$  (resp.  $\bar{\mathcal{F}}_-$ ) to be the set of functions  $f_+$  (resp.  $f_-$ ) belonging to a pair associated with  $(z_{i_0}, v(z_{i_0}))$ . We obtain by construction:

$$|\bar{\mathcal{F}}_+| = |\bar{\mathcal{F}}_-| \geq \frac{|\bar{\mathcal{F}}|}{3M_\gamma n}. \quad (48)$$

Furthermore, according to Lemma 14, the sets  $\bar{\mathcal{F}}_+^{(\eta)}$  and  $\bar{\mathcal{F}}_-^{(\eta)}$  satisfy:

$$\begin{cases} \forall f_+ \in \bar{\mathcal{F}}_+, & f_+^{(\eta)}(z_{i_0}) - b_{i_0} \geq 1 \\ \forall f_- \in \bar{\mathcal{F}}_-, & \max_{k \neq y_{i_0}} f_-^{(\eta)}(z_{i_0}) + b_{i_0} \geq 1 \end{cases} \quad (49)$$

with  $b_{i_0} = v(z_{i_0}) - 1$ . Since  $\bar{\mathcal{F}}_+^{(\eta)} \cup \bar{\mathcal{F}}_-^{(\eta)} \subset \bar{\mathcal{F}}^{(\eta)}$ , obviously, any pair strongly G-shattered by either  $\bar{\mathcal{F}}_+^{(\eta)}$  or  $\bar{\mathcal{F}}_-^{(\eta)}$  is also strongly G-shattered by  $\bar{\mathcal{F}}^{(\eta)}$ . Furthermore, according to (49),  $\bar{\mathcal{F}}^{(\eta)}$  strongly G-shatters the pair  $(\{z_{i_0}\}, b_{i_0})$  which is strongly G-shattered by neither  $\bar{\mathcal{F}}_+^{(\eta)}$  nor  $\bar{\mathcal{F}}_-^{(\eta)}$ . At last, let us consider any pair  $(s'_{\mathcal{Z}^u}, \mathbf{b}'_u)$  strongly G-shattered by both  $\bar{\mathcal{F}}_+^{(\eta)}$  and  $\bar{\mathcal{F}}_-^{(\eta)}$ . Let the pair  $(s''_{\mathcal{Z}^{u+1}}, \mathbf{b}''_{u+1})$  be such that  $s''_{\mathcal{Z}^{u+1}} = s'_{\mathcal{Z}^u} \cup \{z_{i_0}\}$  and the vector  $\mathbf{b}''_{u+1}$  is deduced from  $\mathbf{b}'_u$  by inserting the component  $b_{i_0}$  at the right place. Clearly, neither  $\bar{\mathcal{F}}_+^{(\eta)}$  nor  $\bar{\mathcal{F}}_-^{(\eta)}$  strongly G-shatters  $(s''_{\mathcal{Z}^{u+1}}, \mathbf{b}''_{u+1})$ , simply because they do not strongly G-shatter

the pair  $(\{z_{i_0}\}, b_{i_0})$ . On the contrary, it springs once more from (49) that  $(s''_{\mathcal{Z}^{u+1}}, \mathbf{b}''_{u+1})$  is strongly G-shattered by  $\bar{\mathcal{F}}^{(\eta)}$ . Summarizing, for each pair  $(s'_{\mathcal{Z}^u}, \mathbf{b}'_u)$  strongly G-shattered by both  $\bar{\mathcal{F}}_+^{(\eta)}$  and  $\bar{\mathcal{F}}_-^{(\eta)}$ , we can exhibit by means of an injective mapping a pair  $(s''_{\mathcal{Z}^{u+1}}, \mathbf{b}''_{u+1})$  strongly G-shattered by  $\bar{\mathcal{F}}^{(\eta)}$  but not by  $\bar{\mathcal{F}}_+^{(\eta)}$  or  $\bar{\mathcal{F}}_-^{(\eta)}$ . Collecting all terms, we obtain

$$\begin{aligned} s(\bar{\mathcal{F}}^{(\eta)}) &\geq s(\bar{\mathcal{F}}_+^{(\eta)}) + s(\bar{\mathcal{F}}_-^{(\eta)}) + 1 \\ &\geq \ell(\bar{\mathcal{F}}) - 1, \end{aligned} \tag{50}$$

where the function  $\ell$  returns the number of leaves of the (sub)tree whose root is its argument. Thus, finishing the proof boils down to exhibiting the appropriate lower bound on  $\ell(\bar{\mathcal{F}})$ . To that end, we proceed by induction on the depth of the node. The hypothesis is that

$$\ell(\bar{\mathcal{F}}) \geq |\bar{\mathcal{F}}|^{\frac{1}{\log_2(3M_\gamma n)}}. \tag{51}$$

It is obviously true for the leaves (which are of cardinality 1). Suppose now that it is true for the two sons of an inner node. Then, Inequality (48) gives:

$$\begin{aligned} \ell(\bar{\mathcal{F}}) &= \ell(\bar{\mathcal{F}}_+) + \ell(\bar{\mathcal{F}}_-) \\ &\geq |\bar{\mathcal{F}}_+|^{\frac{1}{\log_2(3M_\gamma n)}} + |\bar{\mathcal{F}}_-|^{\frac{1}{\log_2(3M_\gamma n)}} \\ &\geq 2 \left( \frac{|\bar{\mathcal{F}}|}{3M_\gamma n} \right)^{\frac{1}{\log_2(3M_\gamma n)}} \\ &= |\bar{\mathcal{F}}|^{\frac{1}{\log_2(3M_\gamma n)}}. \end{aligned}$$

The induction hypothesis has been proved. Combining Inequalities (47), (50) and (51) produces by transitivity:

$$|\tilde{\mathcal{F}}|^{\frac{1}{\log_2(3M_\gamma n)}} \leq \Sigma,$$

or equivalently

$$\begin{aligned} |\tilde{\mathcal{F}}| &\leq \Sigma^{\log_2(3M_\gamma n)} \\ &= (3M_\gamma n)^{\log_2(\Sigma)}, \end{aligned}$$

i.e., Inequality (46), the result announced.  $\square$

With Lemma 15 at hand, the proof of Lemma 6 is straightforward.

*Proof.* Let us consider any vector  $\mathbf{z}_n \in \mathcal{Z}^n$  and let  $s_{\mathcal{Z}^n} = \{z_i : 1 \leq i \leq n\}$  be the smallest subset of  $\mathcal{Z}$  containing all the components of  $\mathbf{z}_n$ . Note that its cardinality can be strictly

inferior to  $n$ , in case that  $\mathbf{z}_n$  has two identical components. Setting  $\eta = \frac{\epsilon}{2}$  in (42), one obtains:

$$\mathcal{M}(\epsilon, \mathcal{F}_\gamma, d_{\infty, \mathbf{z}_n}) \leq \mathcal{M}\left(2, \mathcal{F}_\gamma^{(\frac{\epsilon}{2})}, d_{\infty, \mathbf{z}_n}\right).$$

Furthermore, by definition,

$$\mathcal{M}\left(2, \mathcal{F}_\gamma^{(\frac{\epsilon}{2})}, d_{\infty, \mathbf{z}_n}\right) = \mathcal{M}\left(2, \mathcal{F}_\gamma^{(\frac{\epsilon}{2})}\Big|_{s_{\mathcal{Z}^n}}, d_{\infty, \mathbf{z}_n}\right).$$

Let  $\tilde{\mathcal{F}}$  be a subset of  $\mathcal{F}$  such that  $\tilde{\mathcal{F}}_\gamma^{(\frac{\epsilon}{2})}\Big|_{s_{\mathcal{Z}^n}}$  is 2-separated in the metric  $d_{\infty, \mathbf{z}_n}$  and of cardinality  $\mathcal{M}\left(2, \mathcal{F}_\gamma^{(\frac{\epsilon}{2})}\Big|_{s_{\mathcal{Z}^n}}, d_{\infty, \mathbf{z}_n}\right)$  (maximal cardinality).  $\tilde{\mathcal{F}}$  has been built so as to satisfy the hypotheses of Lemma 15. Since  $\left|\tilde{\mathcal{F}}_\gamma^{(\frac{\epsilon}{2})}\Big|_{s_{\mathcal{Z}^n}}\right| = \mathcal{M}\left(2, \mathcal{F}_\gamma^{(\frac{\epsilon}{2})}\Big|_{s_{\mathcal{Z}^n}}, d_{\infty, \mathbf{z}_n}\right)$  and  $\text{S-G-dim}\left(\tilde{\mathcal{F}}^{(\frac{\epsilon}{2})}\right) \leq \text{S-G-dim}\left(\mathcal{F}^{(\frac{\epsilon}{2})}\right)$ , the application of the lemma provides us with:

$$\mathcal{M}\left(2, \mathcal{F}_\gamma^{(\frac{\epsilon}{2})}\Big|_{s_{\mathcal{Z}^n}}, d_{\infty, \mathbf{z}_n}\right) \leq \left(\frac{6\gamma n}{\epsilon}\right)^{\log_2(\Sigma)} \quad (52)$$

where  $\Sigma = \sum_{u=0}^{d_G} \binom{n}{u} \left(\frac{2\gamma}{\epsilon}\right)^u$  with  $d_G$  standing for  $\text{S-G-dim}\left(\mathcal{F}^{(\frac{\epsilon}{2})}\right)$ . According to (45),

$$\text{S-G-dim}\left(\mathcal{F}^{(\frac{\epsilon}{2})}\right) \leq d_G\left(\frac{\epsilon}{4}\right).$$

Since by hypothesis,  $n \geq d_G\left(\frac{\epsilon}{4}\right)$ ,  $\Sigma$  can be bounded from above by replacing in its formula  $d_G$  with  $d_G\left(\frac{\epsilon}{4}\right)$  and resorting to Corollary 3.18 in Mohri et al. (2018), leading to:

$$\begin{aligned} \Sigma &\leq \sum_{u=0}^{d_G\left(\frac{\epsilon}{4}\right)} \binom{n}{u} \left(\frac{2\gamma}{\epsilon}\right)^u \\ &\leq \left(\frac{2\gamma en}{d_G\left(\frac{\epsilon}{4}\right) \epsilon}\right)^{d_G\left(\frac{\epsilon}{4}\right)}, \end{aligned} \quad (53)$$

where the standard convention that the last term takes the value 1 for  $d_G\left(\frac{\epsilon}{4}\right) = 0$  is made. Substituting (53) into (52) and taking the supremum over  $\mathcal{Z}^n$  concludes the proof of (11).  $\square$

### B.3 Margin Graph Dimension - $L_2$ -norm

The sketches of the proofs of the two  $L_2$ -norm combinatorial results, Lemma 7 and Lemma 9, are basically the same. Compared to the sketch of the proof of Lemma 6, they exhibit two major differences. First, the construction of the 2-separating tree is more sophisticated, since it rests on a small deviation principle (in place of the sole pigeonhole principle). Second, one additional step is involved, which implements a probabilistic extraction principle.

This additional step makes the result dimension free. We begin the proof with the formulation of the small deviation principle. This extension of Lemma 5 in Mendelson and Vershynin (2003) is tailored to our needs.

**Lemma 16.** *Let  $T$  be a random variable taking values in  $\llbracket 0; M \rrbracket$  with  $M \geq 2$ . Suppose that  $\text{Var}[T] \geq 9$ . Then there exists either  $(\alpha, \beta) \in \llbracket 1; M-1 \rrbracket \times \left[\frac{1}{M^2}, \frac{1}{2}\right]$  such that*

$$\begin{cases} \mathbb{P}\{T \geq \alpha + 1\} \geq \max\left\{\frac{1}{2}\beta, \frac{1}{M^2}\right\} \\ \mathbb{P}\{T \leq \alpha - 1\} \geq 1 - \beta \end{cases}$$

or  $(\alpha', \beta') \in \llbracket 1; M-1 \rrbracket \times \left[\frac{1}{M^2}, \frac{1}{2}\right]$  such that

$$\begin{cases} \mathbb{P}\{T \geq \alpha' + 1\} \geq 1 - \beta' \\ \mathbb{P}\{T \leq \alpha' - 1\} \geq \max\left\{\frac{1}{2}\beta', \frac{1}{M^2}\right\} \end{cases}.$$

*Proof.* We first note that the hypothesis  $\text{Var}[T] \geq 9$  implies that  $M \geq 6$ . Let  $M_T$  be the smallest median of  $T$  belonging to  $\llbracket 0; M \rrbracket$ . Then, several cases must be distinguished, according to the values of  $M_T$  and  $M - M_T$ . Since they can all be treated the same and the one implying the largest upper bound on the variance, i.e., the one from which springs the hypothesis on the variance, is  $M - M_T \geq 2$  and  $M_T \geq 2$ , we focus on it in the sequel. Let us define the sequences  $(\beta_k)_{k \in \mathbb{N}^*}$  and  $(\beta'_k)_{k \in \mathbb{N}^*}$  as follows:

$$\forall k \in \mathbb{N}^*, \quad \begin{cases} \beta_k = \mathbb{P}\{T \geq M_T + k\} \\ \beta'_k = \mathbb{P}\{T \leq M_T - k\} \end{cases}.$$

Note that by definition of  $M_T$ , both  $\beta_1$  and  $\beta'_1$  are inferior or equal to  $\frac{1}{2}$ . Assume that the conclusion of the lemma fails. We claim that

$$\begin{cases} \forall k \in \llbracket 2; M - M_T \rrbracket, \quad \beta_k \leq \max\left\{\frac{2(M - M_T) - k + 1}{(M - M_T)^3}, \frac{1}{2^k}\right\} \\ \forall k \in \llbracket 2; M_T \rrbracket, \quad \beta'_k \leq \max\left\{\frac{2M_T - k + 1}{M_T^3}, \frac{1}{2^k}\right\} \end{cases}.$$

Indeed, assume that  $\beta_k > \max\left\{\frac{2(M - M_T) - k + 1}{(M - M_T)^3}, \frac{1}{2^k}\right\}$  for some  $k \in \llbracket 2; M - M_T \rrbracket$  and let  $k_0$  be the smallest such index. By construction,  $\beta_{k_0} > \max\left\{\frac{1}{2}\beta_{k_0-1}, \frac{1}{(M - M_T)^2}\right\}$  (even for  $k_0 = 2$  and  $k_0 = M - M_T$ ), so that

$$\begin{cases} \mathbb{P}\{T \geq M_T + k_0\} = \beta_{k_0} > \max\left\{\frac{1}{2}\beta_{k_0-1}, \frac{1}{M^2}\right\} \\ \mathbb{P}\{T \leq M_T + k_0 - 2\} = 1 - \mathbb{P}\{T \geq M_T + k_0 - 1\} = 1 - \beta_{k_0-1} \end{cases}.$$

Since  $\beta_{k_0-1} \geq \beta_{k_0} > \frac{1}{M^2}$  and  $\beta_{k_0-1} \leq \beta_1 \leq \frac{1}{2}$ , so that  $\beta_{k_0-1} \in [\frac{1}{M^2}, \frac{1}{2}]$ , this implies that the conclusion of the lemma would hold with  $\alpha$  being  $M_T + k_0 - 1$  and  $\beta = \beta_{k_0-1}$ , which contradicts the assumption that the conclusion of the lemma fails. The inequality  $\beta'_k \leq \max\left\{\frac{2M_T-k+1}{M_T^3}, \frac{1}{2^k}\right\}$  can be proved in a symmetrical way. As a consequence, upper bounding the maxima by the corresponding sums gives:

$$\begin{aligned}
\text{Var}[T] &= \text{Var}[T - M_T] \\
&\leq \mathbb{E}\left[(T - M_T)^2\right] \\
&= \sum_{t=0}^{+\infty} \mathbb{P}\left\{(T - M_T)^2 > t\right\} \\
&= \sum_{t=0}^{+\infty} \left(\mathbb{P}\left\{T > M_T + \sqrt{t}\right\} + \mathbb{P}\left\{T < M_T - \sqrt{t}\right\}\right) \\
&= \sum_{k=1}^{M-M_T} (2k-1) \beta_k + \sum_{k=1}^{M_T} (2k-1) \beta'_k \\
&< 2 \sum_{k=1}^{+\infty} \frac{2k-1}{2^k} + \sum_{k=2}^{M-M_T} (2k-1) \frac{2(M-M_T)-k+1}{(M-M_T)^3} + \sum_{k=2}^{M_T} (2k-1) \frac{2M_T-k+1}{M_T^3} \\
&\leq 6 + 2 \max_{\Delta \in \mathbb{N} \setminus \{1,2\}} \frac{(8\Delta+11)(\Delta-1)}{6\Delta^2} \\
&< 9.
\end{aligned}$$

This is in contradiction with the hypothesis that  $\text{Var}[T] \geq 9$  and thus concludes the proof.  $\square$

**Proposition 1.** *Let  $\mathcal{T} = \{t_i : 1 \leq i \leq n\}$  be a finite set and  $\mathbf{t}_n = (t_i)_{1 \leq i \leq n}$ . Let  $\mathcal{F}$  be a class of functions from  $\mathcal{T}$  into  $\llbracket 0; M_{\mathcal{F}} \rrbracket$  with  $M_{\mathcal{F}} \geq 2$ . Suppose that  $\mathcal{F}$  is of cardinality at least 2 and  $\epsilon \in \mathbb{R}_+^*$  is such that  $\mathcal{F}$  is  $\epsilon$ -separated in the metric  $d_{2, \mathbf{t}_n}$ . Then there exists  $i_0 \in \llbracket 1; n \rrbracket$  such that*

$$\text{Var}[f(t_{i_0})] \geq \frac{\epsilon^2}{4}.$$

*Proof.* Let us endow  $\mathcal{F}$  with the uniform (counting) measure. Then, the separation assumption on  $\mathcal{F}$  can be used to derive a lower bound on  $\mathbb{E}[d_{2, \mathbf{t}_n}^2(f, f')]$ . Indeed, with probability  $1 - |\mathcal{F}|^{-1}$  we have  $f \neq f'$  and, whenever this event occurs,  $d_{2, \mathbf{t}_n}(f, f') \geq \epsilon$ . As a consequence,

$$\begin{aligned}
\mathbb{E}[d_{2, \mathbf{t}_n}^2(f, f')] &\geq \left(1 - |\mathcal{F}|^{-1}\right) \epsilon^2 \\
&\geq \frac{\epsilon^2}{2}.
\end{aligned}$$

Furthermore,

$$\begin{aligned}\mathbb{E} [d_{2, \mathbf{t}_n}^2 (f, f')] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [(f(t_i) - f'(t_i))^2] \\ &= \frac{2}{n} \sum_{i=1}^n \text{Var} [f(t_i)].\end{aligned}$$

Thus, there exists  $i_0 \in \llbracket 1; n \rrbracket$  such that

$$\text{Var} [f(t_{i_0})] \geq \frac{1}{2} \mathbb{E} [d_{2, \mathbf{t}_n}^2 (f, f')] \geq \frac{\epsilon^2}{4}.$$

□

**Lemma 17.** *Let  $\mathcal{T} = \{t_i : 1 \leq i \leq n\}$  be a finite set and  $\mathbf{t}_n = (t_i)_{1 \leq i \leq n}$ . Let  $\mathcal{F}$  be a class of functions from  $\mathcal{T}$  into  $\llbracket 0; M_{\mathcal{F}} \rrbracket$  with  $M_{\mathcal{F}} \geq 2$ . Suppose that  $\mathcal{F}$  is of cardinality at least 2 and is 6-separated in the metric  $d_{2, \mathbf{t}_n}$ . Then there exist an index  $i_0 \in \llbracket 1; n \rrbracket$  and either  $(\alpha, \beta) \in \llbracket 1; M_{\mathcal{F}} - 1 \rrbracket \times \left[ \frac{1}{M_{\mathcal{F}}^2}, \frac{1}{2} \right]$  such that*

$$\begin{cases} |\{f \in \mathcal{F} : f(t_{i_0}) \geq \alpha + 1\}| \geq \max \left\{ \frac{1}{2}\beta, \frac{1}{M_{\mathcal{F}}^2} \right\} |\mathcal{F}| \\ |\{f \in \mathcal{F} : f(t_{i_0}) \leq \alpha - 1\}| \geq (1 - \beta) |\mathcal{F}| \end{cases}$$

or  $(\alpha', \beta') \in \llbracket 1; M_{\mathcal{F}} - 1 \rrbracket \times \left[ \frac{1}{M_{\mathcal{F}}^2}, \frac{1}{2} \right]$  such that

$$\begin{cases} |\{f \in \mathcal{F} : f(t_{i_0}) \geq \alpha' + 1\}| \geq (1 - \beta') |\mathcal{F}| \\ |\{f \in \mathcal{F} : f(t_{i_0}) \leq \alpha' - 1\}| \geq \max \left\{ \frac{1}{2}\beta', \frac{1}{M_{\mathcal{F}}^2} \right\} |\mathcal{F}| \end{cases}.$$

*Proof.* According to Proposition 1, there exists  $i_0 \in \llbracket 1; n \rrbracket$  such that

$$\text{Var} [f(t_{i_0})] \geq \frac{1}{2} \mathbb{E} [d_{2, \mathbf{t}_n}^2 (f, f')] \geq 9.$$

This implies that the random variable  $f(t_{i_0})$  satisfies the hypotheses of Lemma 16, and the conclusion then springs from the application of this lemma. □

Lemma 17 will be used in the proof of the combinatorial result involving the margin Natarajan dimension: Lemma 9. However, we established it in this section, because its proof can be easily simplified to produce the following variant, appropriate for the margin Graph dimension.

**Lemma 18.** Let  $\mathcal{T} = \{t_i : 1 \leq i \leq n\}$  be a finite set and  $\mathbf{t}_n = (t_i)_{1 \leq i \leq n}$ . Suppose that  $\mathcal{F} \subset \mathbb{Z}^{\mathcal{T}}$  is of cardinality at least 2 and is 5-separated in the metric  $d_{2, \mathbf{t}_n}$ . Then there exist an index  $i_0 \in \llbracket 1; n \rrbracket$  and either  $(\alpha, \beta) \in \mathbb{Z} \times (0, \frac{1}{2}]$  such that

$$\begin{cases} |\{f \in \mathcal{F} : f(t_{i_0}) \geq \alpha + 1\}| \geq \frac{1}{2}\beta |\mathcal{F}| \\ |\{f \in \mathcal{F} : f(t_{i_0}) \leq \alpha - 1\}| \geq (1 - \beta) |\mathcal{F}| \end{cases}$$

or  $(\alpha', \beta') \in \mathbb{Z} \times (0, \frac{1}{2}]$  such that

$$\begin{cases} |\{f \in \mathcal{F} : f(t_{i_0}) \geq \alpha' + 1\}| \geq (1 - \beta') |\mathcal{F}| \\ |\{f \in \mathcal{F} : f(t_{i_0}) \leq \alpha' - 1\}| \geq \frac{1}{2}\beta' |\mathcal{F}| \end{cases}.$$

The following lemma is the basic combinatorial result underlying Lemma 7.

**Lemma 19.** Let  $\mathcal{F}$  be a function class defined as in Definition 8. For every  $s_{\mathcal{Z}^n} = \{z_i = (x_i, y_i) : 1 \leq i \leq n\} \subset \mathcal{Z}$ ,  $\gamma \in (0, 1]$  and  $\eta \in (0, \frac{2}{2}]$ , if  $\mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}}$  is 5-separated in the metric  $d_{2, \mathbf{z}_n}$ , then

$$\left| \mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}} \right| \leq \Sigma^{\frac{3}{2}}, \quad (54)$$

where  $\Sigma = \sum_{u=0}^{d_G} \binom{n}{u} M_\gamma^u$  with  $M_\gamma = \lfloor \frac{2}{\eta} \rfloor$  and  $d_G = S\text{-}G\text{-dim}(\mathcal{F}^{(\eta)})$ .

*Proof.* The principle of the proof is the one of the proof of Lemma 15. Two of the three main formulas still apply: Inequalities (47) and (50). For  $\left| \mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}} \right| \geq 2$ , the incidence of the change of metric is concentrated in the derivation of the 2-separating tree of  $\tilde{\mathcal{F}}$  (or equivalently  $\mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}}$ ), and thus the lower bound on  $\ell(\tilde{\mathcal{F}})$ . Since the inner nodes  $\bar{\mathcal{F}}$  are such that the corresponding sets  $\bar{\mathcal{F}}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}}$  are 5-separated in the metric  $d_{2, \mathbf{z}_n}$ , then according to Lemma 18, for each of these nodes, we can ensure that there exists  $\beta \in (0, \frac{1}{2}]$  such that the two sons  $\bar{\mathcal{F}}_+$  and  $\bar{\mathcal{F}}_-$  of  $\bar{\mathcal{F}}$  verify either  $|\bar{\mathcal{F}}_+| \geq (1 - \beta) |\bar{\mathcal{F}}|$  and  $|\bar{\mathcal{F}}_-| \geq \frac{1}{2}\beta |\bar{\mathcal{F}}|$  or vice versa (in place of (48)). As a consequence, the counterpart of (51) is:

$$\ell(\bar{\mathcal{F}}) \geq |\bar{\mathcal{F}}|^{\frac{2}{3}}. \quad (55)$$

Once more, the proof is an induction on the depth of the node. Inequality (55) is obviously true for the leaves (which are of cardinality 1). Suppose now that it is true for the two sons of an inner node. Then,

$$\begin{aligned} \ell(\bar{\mathcal{F}}) &= \ell(\bar{\mathcal{F}}_+) + \ell(\bar{\mathcal{F}}_-) \\ &\geq \left[ (1 - \beta)^{\frac{2}{3}} + \left(\frac{\beta}{2}\right)^{\frac{2}{3}} \right] |\bar{\mathcal{F}}|^{\frac{2}{3}} \\ &\geq |\bar{\mathcal{F}}|^{\frac{2}{3}}. \end{aligned}$$

Note that the value of the constant,  $\frac{2}{3}$ , can be obtained by maximizing  $K$  over  $(0, 1]$  subject to  $(\frac{1}{2})^K \left[1 + (\frac{1}{2})^K\right] \geq 1$ . Finally, combining Inequalities (47), (50) and (55) produces (54) by transitivity.  $\square$

The following lemma, a slight improvement of Lemma 13 in Mendelson and Vershynin (2003), implements the probabilistic extraction principle.

**Lemma 20.** *Let  $\mathcal{T} = \{t_i : 1 \leq i \leq n\}$  be a finite set,  $\mathbf{t}_n = (t_i)_{1 \leq i \leq n}$  and  $M_{\mathcal{F}} \in \mathbb{R}_+^*$ . Let  $\mathcal{F}$  be a class of functions from  $\mathcal{T}$  into  $[0, M_{\mathcal{F}}]$  with finite cardinality  $|\mathcal{F}| \geq 2$ . Assume that for some  $\epsilon \in (0, M_{\mathcal{F}}]$ ,  $\mathcal{F}$  is  $\epsilon$ -separated with respect to the metric  $d_{2, \mathbf{t}_n}$ , and let*

$$r = \frac{\ln(|\mathcal{F}|)}{K_e \epsilon^4}$$

with

$$K_e = \frac{3}{112M_{\mathcal{F}}^4}.$$

Then, there exists a subvector  $\mathbf{t}'_q$  of  $\mathbf{t}_n$  of size  $q \leq r$  such that  $\mathcal{F}$  is  $\frac{\epsilon}{2}$ -separated with respect to the metric  $d_{2, \mathbf{t}'_q}$ .

*Proof.* This proof uses an abuse of notation that will be repeated in the sequel: the symbol  $\mathbb{P}$  designates different probability measures, some of which implicitly defined. We first note that the statement is trivially true for  $r \geq n$  (it suffices to set  $\mathbf{t}'_q = \mathbf{t}_n$ ). Thus, we proceed under the hypothesis  $r \in [1, n)$ . Let us set  $\mathcal{F} = \{f_j : 1 \leq j \leq |\mathcal{F}|\}$  and  $\mathcal{D}_{\mathcal{F}} = \{f_j - f_{j'} : 1 \leq j < j' \leq |\mathcal{F}|\}$ . The set  $\mathcal{D}_{\mathcal{F}}$  has cardinality  $|\mathcal{D}_{\mathcal{F}}| < \frac{1}{2} |\mathcal{F}|^2$ . Let  $(\epsilon_i)_{1 \leq i \leq n}$  be a sequence of  $n$  independent Bernoulli random variables with common expectation  $\mu = \frac{r}{2n}$ . Then, by application of the  $\epsilon$ -separation property, for every  $\delta_f \in \mathcal{D}_{\mathcal{F}}$ ,

$$\mathbb{P} \left( \frac{1}{n} \sum_{i=1}^n \epsilon_i \delta_f(t_i)^2 < \frac{\epsilon^2 \mu}{2} \right) \leq \mathbb{P} \left( \frac{1}{n} \sum_{i=1}^n (\mu - \epsilon_i) \delta_f(t_i)^2 > \frac{\epsilon^2 \mu}{2} \right). \quad (56)$$

Since by construction, for every  $i \in \llbracket 1; n \rrbracket$ ,  $\mathbb{E} \left[ (\mu - \epsilon_i) \delta_f(t_i)^2 \right] = 0$  and  $|\mu - \epsilon_i| \delta_f(t_i)^2 \leq M_{\mathcal{F}}^2 (1 - \mu) < M_{\mathcal{F}}^2$  with probability one, the right-hand side of (56) can be bounded from above thanks to Bernstein's inequality. Given that

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[ (\mu - \epsilon_i)^2 \delta_f(t_i)^4 \right] \leq M_{\mathcal{F}}^4 \mu (1 - \mu) < M_{\mathcal{F}}^4 \mu,$$

we obtain



$$\begin{aligned}
\mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n \epsilon_i \delta_f(t_i)^2 < \frac{\epsilon^2 \mu}{2}\right) &\leq \mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n (\mu - \epsilon_i) \delta_f(t_i)^2 > \frac{\epsilon^2 \mu}{2}\right) \\
&\leq \exp\left(-\frac{3\mu n \epsilon^4}{4(6M_{\mathcal{F}}^4 + M_{\mathcal{F}}^2 \epsilon^2)}\right) \\
&\leq \exp\left(-\frac{3r \epsilon^4}{56M_{\mathcal{F}}^4}\right) \\
&= |\mathcal{F}|^{-2}.
\end{aligned}$$

Therefore, given the assumption on  $r$ , applying the union bound provides us with:

$$\begin{aligned}
\mathbb{P}\left(\exists \delta_f \in \mathcal{D}_{\mathcal{F}} : \left(\frac{1}{r}\sum_{i=1}^n \epsilon_i \delta_f(t_i)^2\right)^{\frac{1}{2}} < \frac{\epsilon}{2}\right) &= \mathbb{P}\left(\exists \delta_f \in \mathcal{D}_{\mathcal{F}} : \frac{1}{n}\sum_{i=1}^n \epsilon_i \delta_f(t_i)^2 < \frac{\epsilon^2 \mu}{2}\right) \\
&\leq \sum_{\delta_f \in \mathcal{D}_{\mathcal{F}}} \mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n \epsilon_i \delta_f(t_i)^2 < \frac{\epsilon^2 \mu}{2}\right) \\
&\leq |\mathcal{D}_{\mathcal{F}}| \cdot |\mathcal{F}|^{-2} \\
&< \frac{1}{2}. \tag{57}
\end{aligned}$$

Moreover, if  $\mathcal{S}_1$  is the random set  $\{i \in \llbracket 1; n \rrbracket : \epsilon_i = 1\}$ , then by Markov's inequality,

$$\mathbb{P}(|\mathcal{S}_1| > r) = \mathbb{P}\left(\sum_{i=1}^n \epsilon_i > r\right) \leq \frac{1}{2}. \tag{58}$$

Combining (57) and (58) by means of the union bound provides us with

$$\mathbb{P}\left\{\left(\exists \delta_f \in \mathcal{D}_{\mathcal{F}} : \left(\frac{1}{r}\sum_{i=1}^n \epsilon_i \delta_f(t_i)^2\right)^{\frac{1}{2}} < \frac{\epsilon}{2}\right) \text{ or } (|\mathcal{S}_1| > r)\right\} < 1$$

which implies that

$$\mathbb{P}\left\{\left(\forall \delta_f \in \mathcal{D}_{\mathcal{F}} : \|\delta_f\|_{L_2(\mu_{(t_i)_{i \in \mathcal{S}_1})})} \geq \frac{\epsilon}{2}\right) \text{ and } (|\mathcal{S}_1| \leq r)\right\} > 0.$$

This translates into the fact that there exists a subvector  $\mathbf{t}'_q$  of  $\mathbf{t}_n$  of size  $q \leq r$  such that the class  $\mathcal{F}$  is  $\frac{\epsilon}{2}$ -separated with respect to the metric  $d_{2, \mathbf{t}'_q}$ , i.e., our claim.  $\square$

The proof of Lemma 7 is the following one.

*Proof.* Let us consider any vector  $\mathbf{z}_n \in \mathcal{Z}^n$  and let  $s_{\mathcal{Z}^n} = \{z_i : 1 \leq i \leq n\}$  be the smallest subset of  $\mathcal{Z}$  containing all the components of  $\mathbf{z}_n$ . Note that its cardinality can be strictly inferior to  $n$ , in case that  $\mathbf{z}_n$  has two identical components. By definition,

$$\mathcal{M}(\epsilon, \mathcal{F}_\gamma, d_{2, \mathbf{z}_n}) = \mathcal{M}\left(\epsilon, \mathcal{F}_\gamma|_{s_{\mathcal{Z}^n}}, d_{2, \mathbf{z}_n}\right).$$

Let  $\tilde{\mathcal{F}}$  be a subset of  $\mathcal{F}$  of cardinality  $\mathcal{M}(\epsilon, \mathcal{F}_\gamma, d_{2, \mathbf{z}_n})$  such that  $\tilde{\mathcal{F}}_\gamma|_{s_{\mathcal{Z}^n}}$  is  $\epsilon$ -separated with respect to the metric  $d_{2, \mathbf{z}_n}$  and in bijection with  $\tilde{\mathcal{F}}$ . By construction,  $\tilde{\mathcal{F}}_\gamma|_{s_{\mathcal{Z}^n}}$  satisfies the hypotheses of Lemma 20, with  $M_{\mathcal{F}} = \gamma$ . Consequently, applying the lemma establishes that there exists a subvector  $\mathbf{z}'_q$  of  $\mathbf{z}_n$  of size

$$q \leq \frac{\ln\left(|\tilde{\mathcal{F}}|\right)}{K_\gamma \epsilon^4} \quad (59)$$

with  $K_\gamma = \frac{3}{112\gamma^4}$  such that the class  $\tilde{\mathcal{F}}_\gamma|_{s_{\mathcal{Z}^n}}$  is also  $\frac{\epsilon}{2}$ -separated with respect to the metric  $d_{2, \mathbf{z}'_q}$ . As a consequence, denoting  $s_{\mathcal{Z}^q} = \{z'_i = (x'_i, y'_i) : 1 \leq i \leq q\}$  ( $|s_{\mathcal{Z}^q}| \leq q$ ), it appears that  $|\tilde{\mathcal{F}}_\gamma|_{s_{\mathcal{Z}^n}}| = |\tilde{\mathcal{F}}_\gamma|_{s_{\mathcal{Z}^q}}|$  (and thus  $|\tilde{\mathcal{F}}_\gamma|_{s_{\mathcal{Z}^q}}| = |\tilde{\mathcal{F}}|$ ). Applying (41) to the function class  $\tilde{\mathcal{F}}_\gamma|_{s_{\mathcal{Z}^q}}$  with  $N = 5$  and the corresponding largest possible value for  $\eta$ ,  $\frac{\epsilon}{12}$ , it appears that  $\tilde{\mathcal{F}}_\gamma^{(\frac{\epsilon}{12})}|_{s_{\mathcal{Z}^q}}$  is a set of cardinality  $|\tilde{\mathcal{F}}|$  which is 5-separated with respect to the metric  $d_{2, \mathbf{z}'_q}$ . Thus, Lemma 19 can be applied to  $\tilde{\mathcal{F}}$  as follows:

$$\begin{aligned} |\tilde{\mathcal{F}}| &= \left| \tilde{\mathcal{F}}_\gamma^{(\frac{\epsilon}{12})}|_{s_{\mathcal{Z}^q}} \right| \\ &\leq \left( \sum_{u=0}^{d_G} \binom{q}{u} \left(\frac{12\gamma}{\epsilon}\right)^u \right)^{\frac{3}{2}} \\ &\leq \left( \frac{12\gamma e q}{d_G \epsilon} \right)^{\frac{3d_G}{2}}, \end{aligned} \quad (60)$$

where  $d_G = \text{S-G-dim}\left(\tilde{\mathcal{F}}^{(\frac{\epsilon}{12})}\right)$ . A substitution of the upper bound on  $q$  provided by (59) into (60) gives:

$$|\tilde{\mathcal{F}}| \leq \left( K \left(\frac{\gamma}{\epsilon}\right)^5 \frac{\ln\left(|\tilde{\mathcal{F}}|\right)}{d_G} \right)^{\frac{3}{2} d_G}$$

with  $K = 448e$ . In order to upper bound  $\ln\left(|\tilde{\mathcal{F}}|^{\frac{1}{d_G}}\right)$ , we resort once more to (40), this time with  $u_0 = 1$ . Thus,

$$\left|\tilde{\mathcal{F}}\right|^{\frac{1}{d_G}} \leq \ln^{\frac{3}{2}}\left(\left|\tilde{\mathcal{F}}\right|^{\frac{1}{d_G}}\right) \left(K \left(\frac{\gamma}{\epsilon}\right)^5\right)^{\frac{3}{2}}$$

and  $|\tilde{\mathcal{F}}| = \mathcal{M}(\epsilon, \mathcal{F}_\gamma, d_{2, \mathbf{z}_n})$  imply that

$$\begin{aligned} \mathcal{M}(\epsilon, \mathcal{F}_\gamma, d_{2, \mathbf{z}_n}) &\leq \left(2K \left(\frac{\gamma}{\epsilon}\right)^5\right)^{\frac{12}{5}d_G} \\ &\leq \left(\frac{5\gamma}{\epsilon}\right)^{12d_G}. \end{aligned} \quad (61)$$

Since  $\tilde{\mathcal{F}}^{(\frac{\epsilon}{12})} \subset \mathcal{F}^{(\frac{\epsilon}{12})}$ , by application of Formula (45),

$$\begin{aligned} \text{S-G-dim} \left(\tilde{\mathcal{F}}^{(\frac{\epsilon}{12})}\right) &\leq \text{S-G-dim} \left(\mathcal{F}^{(\frac{\epsilon}{12})}\right) \\ &\leq d_G \left(\frac{\epsilon}{24}\right). \end{aligned} \quad (62)$$

By substitution of (62) into (61), we obtain that for every vector  $\mathbf{z}_n \in \mathcal{Z}^n$ ,

$$\mathcal{M}(\epsilon, \mathcal{F}_\gamma, d_{2, \mathbf{z}_n}) \leq \left(\frac{5\gamma}{\epsilon}\right)^{12d_G \left(\frac{\epsilon}{24}\right)}. \quad (63)$$

At last, (63) implies (12) since its right-hand side does not depend on  $\mathbf{z}_n$ .  $\square$

#### B.4 Margin Natarajan Dimension - Uniform Convergence Norm

The proof of Lemma 8 is essentially that of Lemma 6, with the main differences being concentrated in the basic combinatorial result (the counterpart of Lemma 15). Thus, we only highlight these differences. The first one is that the number  $s(\bar{\mathcal{F}}^{(\eta)})$  of pairs  $(s'_{\mathcal{Z}^u}, \mathbf{b}'_u)$  strongly G-shattered by  $\bar{\mathcal{F}}^{(\eta)}$  is replaced with the number  $s'(\bar{\mathcal{F}}^{(\eta)})$  of triplets  $(s'_{\mathcal{Z}^u}, \mathbf{b}'_u, \mathbf{c}'_u)$  strongly N-shattered by  $\bar{\mathcal{F}}^{(\eta)}$ . Let  $d_N = \text{S-N-dim}(\mathcal{F}^{(\eta)})$ . Once more, under the hypothesis  $d_N \geq 1$ , combinatorics provides us with:

$$s'(\mathcal{F}^{(\eta)}) \leq \sum_{u=1}^{d_N} \binom{n}{u} M_\gamma^u (C-1)^u = \Sigma' - 1. \quad (64)$$

In order to obtain the counterpart of (50), i.e.,

$$\begin{aligned} s'(\bar{\mathcal{F}}^{(\eta)}) &\geq s'(\bar{\mathcal{F}}_+^{(\eta)}) + s'(\bar{\mathcal{F}}_-^{(\eta)}) + 1 \\ &\geq \ell(\bar{\mathcal{F}}) - 1, \end{aligned} \quad (65)$$

(49) must be replaced with

$$\begin{cases} \forall f_+ \in \bar{\mathcal{F}}_+, & f_+^{(\eta)}(z_{i_0}) - b_{i_0} \geq 1 \\ \forall f_- \in \bar{\mathcal{F}}_-, & f_-^{(\eta)}(x_{i_0}, c_{i_0}) + b_{i_0} \geq 1 \end{cases}. \quad (66)$$

This calls for an additional application of the pigeonhole principle in the derivation of the class  $\bar{\mathcal{F}}_-$ , so that the lower bound on  $|\bar{\mathcal{F}}_-|$  provided by (48) is replaced with

$$|\bar{\mathcal{F}}_-| \geq \frac{|\bar{\mathcal{F}}|}{3M_\gamma(C-1)n}.$$

This implies that the counterpart of (51) is

$$\ell(\bar{\mathcal{F}}) \geq |\bar{\mathcal{F}}|^{\frac{1}{\log_2(3M_\gamma\sqrt{C-1}n)}}. \quad (67)$$

Once more, it is proved by induction on the depth of the node. Inequality (67) is obviously true for the leaves (which are of cardinality 1). Suppose now that it is true for the two sons of an inner node. Then,

$$\begin{aligned} \ell(\bar{\mathcal{F}}) &= \ell(\bar{\mathcal{F}}_+) + \ell(\bar{\mathcal{F}}_-) \\ &\geq \left( \frac{|\bar{\mathcal{F}}|}{3M_\gamma n} \right)^{\frac{1}{\log_2(3M_\gamma\sqrt{C-1}n)}} + \left( \frac{|\bar{\mathcal{F}}|}{3M_\gamma(C-1)n} \right)^{\frac{1}{\log_2(3M_\gamma\sqrt{C-1}n)}} \\ &= \frac{1}{2} \left( \left( \sqrt{C-1} \right)^{\frac{1}{\log_2(3M_\gamma\sqrt{C-1}n)}} + \left( \sqrt{C-1} \right)^{-\frac{1}{\log_2(3M_\gamma\sqrt{C-1}n)}} \right) |\bar{\mathcal{F}}|^{\frac{1}{\log_2(3M_\gamma\sqrt{C-1}n)}} \\ &\geq \frac{1}{2} \min_{t \in \mathbb{R}_+^*} \left( t + \frac{1}{t} \right) |\bar{\mathcal{F}}|^{\frac{1}{\log_2(3M_\gamma\sqrt{C-1}n)}} \\ &= |\bar{\mathcal{F}}|^{\frac{1}{\log_2(3M_\gamma\sqrt{C-1}n)}}. \end{aligned}$$

Combining Inequalities (64), (65) and (67), the counterpart of Inequality (46) is

$$\left| \mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}} \right| \leq \left( 3M_\gamma\sqrt{C-1}n \right)^{\log_2(\Sigma')}.$$

## B.5 Margin Natarajan Dimension - $L_2$ -Norm

The main difference between the proof of Lemma 9 and the proof of Lemma 7 is located in the small deviation principle (Lemma 17 replaces Lemma 18). Since the consequences of this change appear in the derivation of the basic combinatorial result, we provide this latter result with its full proof.

**Lemma 21.** *Let  $\mathcal{F}$  be a function class defined as in Definition 8. For every  $s_{\mathcal{Z}^n} = \{z_i = (x_i, y_i) : 1 \leq i \leq n\} \subset \mathcal{Z}$ ,  $\gamma \in (0, 1]$  and  $\eta \in (0, \frac{\gamma}{2}]$ , if  $\mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}}$  is 6-separated in the metric  $d_{2, \mathbf{z}_n}$ , then*

$$\left| \mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}} \right| \leq (\Sigma')^{\log_2 \left( \frac{M_\gamma^2}{\sqrt{M_\gamma^2 - 2}} \sqrt{C-1} \right)} \leq (\Sigma')^{\frac{1}{2} \log_2(2M_\gamma^2(C-1))} \quad (68)$$

where  $\Sigma' = \sum_{u=0}^{d_N} \binom{n}{u} M_\gamma^u (C-1)^u$  with  $M_\gamma = \left\lfloor \frac{\gamma}{\eta} \right\rfloor$  and  $d_N = S\text{-}N\text{-dim}(\mathcal{F}^{(\eta)})$ .

*Proof.* Inequality (68) is trivially true for  $\left| \mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}} \right| = 1$ . Indeed, the minimal value of its right-hand side, corresponding to  $d_N = 0$ , is 1. Thus, the rest of the proof makes use of the restriction  $\left| \mathcal{F}_\gamma^{(\eta)} \Big|_{s_{\mathcal{Z}^n}} \right| \geq 2$ . A direct consequence is that according to Lemma 14,  $d_N \geq 1$ . A subset of  $s_{\mathcal{Z}^n}$  of cardinality  $u \in \llbracket 1; n \rrbracket$  is denoted by  $s'_{\mathcal{Z}^u} = \{z'_i : 1 \leq i \leq u\}$ , with the convention

$$\forall (i, j) : 1 \leq i < j \leq u, (z'_i, z'_j) = (z_v, z_w) \implies 1 \leq v < w \leq n.$$

For every subset  $\bar{\mathcal{F}}$  of  $\mathcal{F}$ , denote by  $s'(\bar{\mathcal{F}}^{(\eta)})$  the number of triplets  $(s'_{\mathcal{Z}^u}, \mathbf{b}'_u, \mathbf{c}'_u)$  with  $s'_{\mathcal{Z}^u} \subset s_{\mathcal{Z}^n}$ ,  $\mathbf{b}'_u \in \llbracket 1; M_\gamma - 1 \rrbracket^u$  and  $\mathbf{c}'_u \in \mathcal{Y}^u$  (with  $\forall i \in \llbracket 1; u \rrbracket, c'_i \neq y'_i$ ) strongly N-shattered by  $\bar{\mathcal{F}}^{(\eta)}$  (the convention above has been introduced to avoid handling duplicates). Since  $d_N \geq 1$ , Inequality (64) provides us once more with an upper bound on  $s'(\mathcal{F}^{(\eta)})$ . In order to derive a lower bound on the same quantity, we also build a 2-separating tree of  $\tilde{\mathcal{F}}$ . Let  $\bar{\mathcal{F}}$  be one of its nodes such that  $|\bar{\mathcal{F}}| \geq 2$  (inner node). Its two sons,  $\bar{\mathcal{F}}_+$  and  $\bar{\mathcal{F}}_-$ , are built by application of Lemma 17 and the pigeonhole principle. According to Lemma 17, we can ensure that there exist an index  $i_0 \in \llbracket 1; n \rrbracket$ ,  $(\alpha, \beta) \in \llbracket 1; M_\gamma - 1 \rrbracket \times \left[ \frac{1}{M_\gamma^2}, \frac{1}{2} \right]$  and two subsets  $\hat{\mathcal{F}}_+$  and  $\hat{\mathcal{F}}_-$  of  $\bar{\mathcal{F}}$  verifying either  $|\hat{\mathcal{F}}_+| \geq (1 - \beta) |\bar{\mathcal{F}}|$  and  $|\hat{\mathcal{F}}_-| \geq \max \left\{ \frac{1}{2}\beta, \frac{1}{M_\gamma^2} \right\} |\bar{\mathcal{F}}|$  or vice versa, such that

$$\begin{cases} \forall f_+ \in \bar{\mathcal{F}}_+, f_{+, \gamma}^{(\eta)}(z_{i_0}) \geq \alpha + 1 \\ \forall f_- \in \hat{\mathcal{F}}_-, f_{-, \gamma}^{(\eta)}(z_{i_0}) \leq \alpha - 1 \end{cases}.$$

Setting  $b_{i_0} = \alpha$ , it springs from Lemma 14 that

$$\begin{cases} \forall f_+ \in \bar{\mathcal{F}}_+, f_+^{(\eta)}(z_{i_0}) - b_{i_0} \geq 1 \\ \forall f_- \in \hat{\mathcal{F}}_-, \max_{k \neq y_{i_0}} f_-^{(\eta)}(x_{i_0}, k) + b_{i_0} \geq 1 \end{cases},$$

i.e., (49) is obtained with  $\bar{\mathcal{F}}_-$  replaced with  $\hat{\mathcal{F}}_-$ . There comes the application of the pigeonhole principle, to obtain Formula (66). The derivation of the corresponding function classes is as follows. There exists  $c_{i_0} \in \mathcal{Y} \setminus \{y_{i_0}\}$  such that among the functions  $f_-$  in  $\hat{\mathcal{F}}_-$ , at least  $\left\lfloor \frac{|\hat{\mathcal{F}}_-|}{C-1} \right\rfloor$  of them satisfy  $c_{i_0} \in \operatorname{argmax}_{k \neq y_{i_0}} f_-(x_{i_0}, k)$ . We choose  $\bar{\mathcal{F}}_-$  to be any such subset of  $\hat{\mathcal{F}}_-$ . With Formula (66) at hand, Inequality (65) is also available. Thus, finishing the proof still boils down to deriving a lower bound on  $\ell(\bar{\mathcal{F}})$ . The originality rests on the fact that two cases must be considered, to take into account the two sources of asymmetry between the cardinalities of  $\bar{\mathcal{F}}_+$  and  $\bar{\mathcal{F}}_-$ . Indeed, we have either  $|\bar{\mathcal{F}}_+| \geq \max \left\{ \frac{1}{2}\beta, \frac{1}{M_\gamma^2} \right\} |\bar{\mathcal{F}}|$  and  $|\bar{\mathcal{F}}_-| \geq \frac{1-\beta}{C-1} |\bar{\mathcal{F}}|$  or  $|\bar{\mathcal{F}}_+| \geq (1 - \beta) |\bar{\mathcal{F}}|$  and  $|\bar{\mathcal{F}}_-| \geq \frac{1}{C-1} \max \left\{ \frac{1}{2}\beta, \frac{1}{M_\gamma^2} \right\} |\bar{\mathcal{F}}|$ . The

induction hypothesis is this time:

$$\ell(\bar{\mathcal{F}}) \geq |\bar{\mathcal{F}}|^{\frac{1}{\log_2\left(\frac{M_\gamma^2}{\sqrt{M_\gamma^2-2}}\sqrt{C-1}\right)}}. \quad (69)$$

Once more, it is obviously true for the leaves. We prove it for the first case (the other one is treated in the same way). The computation makes use of the following analytical result:

$$\forall K \in (0, 1], \quad \operatorname{argmin}_{\beta \in \left[\frac{1}{M_\gamma^2}, \frac{1}{2}\right]} \left\{ \left(\frac{1-\beta}{C-1}\right)^K + \left(\max\left\{\frac{\beta}{2}, \frac{1}{M_\gamma^2}\right\}\right)^K \right\} = \frac{2}{M_\gamma^2}.$$

Then,

$$\begin{aligned} \ell(\bar{\mathcal{F}}) &= \ell(\bar{\mathcal{F}}_+) + \ell(\bar{\mathcal{F}}_-) \\ &\geq \left( \left(\frac{1}{M_\gamma^2}\right)^{\log_2\left(\frac{M_\gamma^2}{\sqrt{M_\gamma^2-2}}\sqrt{C-1}\right)} + \left(\frac{M_\gamma^2-2}{M_\gamma^2(C-1)}\right)^{\log_2\left(\frac{M_\gamma^2}{\sqrt{M_\gamma^2-2}}\sqrt{C-1}\right)} \right) |\bar{\mathcal{F}}|^{\log_2\left(\frac{M_\gamma^2}{\sqrt{M_\gamma^2-2}}\sqrt{C-1}\right)} \\ &= \frac{1}{2} \left( \left(\sqrt{\frac{C-1}{M_\gamma^2-2}}\right)^{\log_2\left(\frac{M_\gamma^2}{\sqrt{M_\gamma^2-2}}\sqrt{C-1}\right)} + \left(\sqrt{\frac{M_\gamma^2-2}{C-1}}\right)^{\log_2\left(\frac{M_\gamma^2}{\sqrt{M_\gamma^2-2}}\sqrt{C-1}\right)} \right) \\ &\quad \times |\bar{\mathcal{F}}|^{\log_2\left(\frac{M_\gamma^2}{\sqrt{M_\gamma^2-2}}\sqrt{C-1}\right)} \\ &\geq \frac{1}{2} \min_{t \in \mathbb{R}_+^*} \left( t + \frac{1}{t} \right) |\bar{\mathcal{F}}|^{\log_2\left(\frac{M_\gamma^2}{\sqrt{M_\gamma^2-2}}\sqrt{C-1}\right)} \\ &= |\bar{\mathcal{F}}|^{\log_2\left(\frac{M_\gamma^2}{\sqrt{M_\gamma^2-2}}\sqrt{C-1}\right)}. \end{aligned}$$

Combining Inequalities (64), (65) and (69) produces by transitivity:

$$|\tilde{\mathcal{F}}| \leq (\Sigma')^{\log_2\left(\frac{M_\gamma^2}{\sqrt{M_\gamma^2-2}}\sqrt{C-1}\right)}.$$

The left-hand side inequality of Formula (68) has been established. To conclude the proof, it suffices to remember that  $M_\gamma \geq 2$ .  $\square$

The proof of Lemma 9 is the following one.

*Proof.* The beginning of the proof is identical to the beginning of the proof of Lemma 7 up to the use of Formula (41), which is now done with  $N = 6$  (instead of  $N = 5$ ). The reason

for this change is to satisfy the hypotheses of the basic combinatorial result, Lemma 21, which replaces Lemma 19.

$$\begin{aligned}
|\tilde{\mathcal{F}}| &= \left| \tilde{\mathcal{F}}_\gamma^{\left(\frac{\epsilon}{14}\right)} \Big|_{s_{Z^q}} \right| \\
&\leq \left( \sum_{u=0}^{d_N} \binom{q}{u} \left(\frac{14\gamma}{\epsilon}\right)^u (C-1)^u \right)^{\frac{1}{2} \log_2 \left( 2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) \right)} \\
&\leq \left( \frac{14\gamma(C-1)eq}{d_N \epsilon} \right)^{\frac{1}{2} \log_2 \left( 2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) \right) d_N}, \tag{70}
\end{aligned}$$

where  $d_N = \text{S-N-dim} \left( \tilde{\mathcal{F}}^{\left(\frac{\epsilon}{14}\right)} \right)$ . A substitution of the upper bound on  $q$  provided by (59) into (70) gives:

$$|\tilde{\mathcal{F}}| \leq \left( K(C-1) \left(\frac{\gamma}{\epsilon}\right)^5 \frac{\ln \left( |\tilde{\mathcal{F}}| \right)}{d_N} \right)^{\frac{1}{2} \log_2 \left( 2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) \right) d_N}$$

with  $K = \frac{1568}{3}e$ . In order to upper bound  $\ln \left( |\tilde{\mathcal{F}}|^{\frac{1}{d_N}} \right)$ , we resort once more to (40), with  $u_0 = \frac{1}{4} \log_2 \left( 2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) \right)$ . Thus,

$$\ln \left( |\tilde{\mathcal{F}}|^{\frac{1}{d_N}} \right) \leq \frac{1}{2} \log_2 \left( 2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) \right) |\tilde{\mathcal{F}}|^{\frac{1}{\log_2 \left( 2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) \right) d_N}}$$

implies that

$$|\tilde{\mathcal{F}}| \leq \left( \frac{1}{2} \log_2 \left( 2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) \right) K(C-1) \left(\frac{\gamma}{\epsilon}\right)^5 \right)^{\log_2 \left( 2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) \right) d_N}.$$

Since

$$\begin{aligned}
2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) > 16 &\implies \log_2 \left( 2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) \right) < \left( 2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) \right)^{\frac{1}{2}} \\
&\implies \log_2 \left( 2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) \right) < \left( \frac{K}{2} (C-1) \left(\frac{\gamma}{\epsilon}\right)^5 \right)^{\frac{1}{2}},
\end{aligned}$$

the upper bound on  $|\tilde{\mathcal{F}}|$  simplifies into

$$\mathcal{M}(\epsilon, \mathcal{F}_\gamma, d_{2, \mathbf{z}_n}) = |\tilde{\mathcal{F}}| \leq \left( (C-1) \left(\frac{4\gamma}{\epsilon}\right)^5 \right)^{\frac{3}{2} \log_2 \left( 2 \left(\frac{14\gamma}{\epsilon}\right)^2 (C-1) \right) d_N}. \tag{71}$$

Since  $\tilde{\mathcal{F}}^{(\frac{\epsilon}{14})} \subset \mathcal{F}^{(\frac{\epsilon}{14})}$ , by application of Formula (45),

$$\begin{aligned} \text{S-N-dim} \left( \tilde{\mathcal{F}}^{(\frac{\epsilon}{14})} \right) &\leq \text{S-N-dim} \left( \mathcal{F}^{(\frac{\epsilon}{14})} \right) \\ &\leq d_N \left( \frac{\epsilon}{28} \right). \end{aligned} \quad (72)$$

A substitution of (72) into (71) produces an upper bound on  $\mathcal{M}(\epsilon, \mathcal{F}_\gamma, d_{2, \mathbf{z}_n})$  which does not depend on  $\mathbf{z}_n$ , thus concluding the proof.  $\square$

## C Proofs of the Structural Results

This appendix gathers the proofs of the upper bounds on  $\gamma$ -N-dim( $\rho_G$ ). The proof of the first of them, Lemma 10, makes use of three partial results which are now stated.

### C.1 Technical Lemmas

**Proposition 2.** *Let  $\mathcal{F}$  be a function class defined as in Definition 8. Suppose that for  $\gamma \in \mathbb{R}_+^*$ , the subset  $\bar{\mathcal{F}}$  of  $\mathcal{F}$   $\gamma$ -N-shatters the triplet  $(\{(x_i, y_i) : 1 \leq i \leq n\}, \mathbf{b}_n, \mathbf{c}_n)$ . Then  $\bar{\mathcal{F}}$  also  $\gamma$ -N-shatters another triplet,  $(\{(x_i, y'_i) : 1 \leq i \leq n\}, \mathbf{b}'_n, \mathbf{c}'_n)$ , derived from the first one as follows:*

$$\forall i \in \llbracket 1; n \rrbracket, \begin{cases} \text{if } y_i < c_i, y'_i = y_i, b'_i = b_i, c'_i = c_i \\ \text{if } y_i > c_i, y'_i = c_i, b'_i = -b_i, c'_i = y_i \end{cases}.$$

As a consequence, the derivation of an upper bound on  $\gamma$ -N-dim( $\mathcal{F}$ ) can make use of a stronger hypothesis on  $(\mathbf{y}_n, \mathbf{c}_n)$ :  $\forall i \in \llbracket 1; n \rrbracket, y_i < c_i$ , provided that the hypothesis of non-negativity of the biases  $b_i$  is relaxed.

*Proof.* Without loss of generality, we assume that  $\bar{\mathcal{F}}$  is of minimal cardinality  $2^n$  and set accordingly  $\bar{\mathcal{F}} = \{f_{\mathbf{s}_n} : \mathbf{s}_n \in \{-1, 1\}^n\}$ . Consider the following bijection on  $\{-1, 1\}^n$ :

$$\begin{aligned} B : \{-1, 1\}^n &\longrightarrow \{-1, 1\}^n \\ \mathbf{s}_n &\mapsto \mathbf{s}'_n \\ \forall i \in \llbracket 1; n \rrbracket, &\begin{cases} \text{if } y_i < c_i, s'_i = s_i \\ \text{if } y_i > c_i, s'_i = -s_i \end{cases}. \end{aligned}$$

Then,

$$\forall i \in \llbracket 1; n \rrbracket, \begin{cases} \text{if } s'_i = 1, f_{\mathbf{s}'_n}(x_i, y_i) - b_i \geq \gamma \\ \text{if } s'_i = -1, f_{\mathbf{s}'_n}(x_i, c_i) + b_i \geq \gamma \end{cases} \implies \forall i \in \llbracket 1; n \rrbracket, \begin{cases} \text{if } s_i = 1, f_{\mathbf{s}'_n}(x_i, y'_i) - b'_i \geq \gamma \\ \text{if } s_i = -1, f_{\mathbf{s}'_n}(x_i, c'_i) + b'_i \geq \gamma \end{cases}.$$



By definition, the triplet  $(s'_{\mathcal{Z}^n}, \mathbf{b}'_n, \mathbf{c}'_n)$  is  $\gamma$ -N-shattered by  $\bar{\mathcal{F}}$ , which concludes the proof.  $\square$

Proposition 3 is an extension of Proposition 1.4 in Talagrand (2003) holding for the  $L_p$ -norms with  $p \in \mathbb{N} \setminus \{0, 1\}$  (instead of simply  $p = 2$ ), that explicits the value of the constant.

**Proposition 3.** *Let  $\mathcal{F}$  be a class of real-valued functions on  $\mathcal{T}$ . For every  $\gamma \in \mathbb{R}_+^*$  satisfying  $\gamma\text{-dim}(\mathcal{F}) > 0$ ,  $n \in \llbracket 1; \gamma\text{-dim}(\mathcal{F}) \rrbracket$  and  $p \in \mathbb{N} \setminus \{0, 1\}$ ,*

$$n \leq K_p \log_2 (\mathcal{M}_p(\gamma, \mathcal{F}, n))$$

$$\text{with } K_p = \left( \frac{2^p}{2^{p-1}-1} \right)^2.$$

*Proof.* Suppose that for  $\gamma \in \mathbb{R}_+^*$ , the subset  $s_{\mathcal{T}^n} = \{t_i : 1 \leq i \leq n\}$  of  $\mathcal{T}$  is  $\gamma$ -shattered by  $\mathcal{F}$  and  $\mathbf{b}_n = (b_i)_{1 \leq i \leq n} \in \mathbb{R}^n$  is a witness to this shattering. By definition, there exists a subset  $\bar{\mathcal{F}} = \{f_{\mathbf{s}_n} : \mathbf{s}_n \in \{-1, 1\}^n\}$  of  $\mathcal{F}$  satisfying

$$\forall \mathbf{s}_n \in \{-1, 1\}^n, \forall i \in \llbracket 1; n \rrbracket, s_i(f_{\mathbf{s}_n}(t_i) - b_i) \geq \gamma. \quad (73)$$

Let  $\mathbf{t}_n = (t_i)_{1 \leq i \leq n}$ . To prove the proposition, it suffices to establish that

$$n \leq K_p \log_2 (\mathcal{M}(\gamma, \bar{\mathcal{F}}, d_{p, \mathbf{t}_n})). \quad (74)$$

For  $(\mathbf{s}_n, \mathbf{s}'_n) \in (\{-1, 1\}^n)^2$ , let  $\mathcal{S}(\mathbf{s}_n, \mathbf{s}'_n)$  be the subset of  $\llbracket 1; n \rrbracket$  defined by:

$$\mathcal{S}(\mathbf{s}_n, \mathbf{s}'_n) = \{i \in \llbracket 1; n \rrbracket : s_i \neq s'_i\}.$$

Then, making use of (73), we obtain that

$$\begin{aligned} d_{p, \mathbf{t}_n}(f_{\mathbf{s}_n}, f_{\mathbf{s}'_n}) &\geq \left( \frac{1}{n} |\mathcal{S}(\mathbf{s}_n, \mathbf{s}'_n)| (2\gamma)^p \right)^{\frac{1}{p}} \\ &= 2\gamma \left( \frac{d_H(\mathbf{s}_n, \mathbf{s}'_n)}{n} \right)^{\frac{1}{p}}, \end{aligned}$$

where  $d_H$  stands for the Hamming distance. Thus, a sufficient condition for  $d_{p, \mathbf{t}_n}(f_{\mathbf{s}_n}, f_{\mathbf{s}'_n}) \geq \gamma$  is  $d_H(\mathbf{s}_n, \mathbf{s}'_n) \geq \lceil (\frac{1}{2})^p n \rceil$ . As a consequence, to prove (74), it suffices to establish that there is a subset of the set of vertices of the hypercube  $Q_n$  of cardinality  $\lceil 2^{\frac{n}{K_p}} \rceil$  which is  $\lceil (\frac{1}{2})^p n \rceil$ -separated with respect to the Hamming distance (the separation is well-defined since  $\lceil 2^{\frac{n}{K_p}} \rceil \geq 2$ ). To that end, a probabilistic approach similar to that of the proof of Lemma 20 is implemented. For  $q \in \llbracket 2; 2^n \rrbracket$ , let  $\epsilon_{q, n} = (\epsilon_{j, i})_{1 \leq j \leq q, 1 \leq i \leq n}$  be a Bernoulli

random matrix (its entries  $\epsilon_{j,i}$  are independent Bernoulli random variables with common expectation  $\frac{1}{2}$ ). Then, by application of the union bound,

$$\begin{aligned} & \mathbb{P} \left( \exists (j, j') \in \llbracket 1; q \rrbracket^2 : 1 \leq j < j' \leq q \text{ and } \sum_{i=1}^n \mathbf{1}_{\{\epsilon_{j,i} \neq \epsilon_{j',i}\}} < \left(\frac{1}{2}\right)^p n \right) \\ & \leq \binom{q}{2} \mathbb{P} \left( \sum_{i=1}^n \epsilon_i > n \left(1 - \left(\frac{1}{2}\right)^p\right) \right), \end{aligned}$$

where  $(\epsilon_i)_{1 \leq i \leq n}$  is a Bernoulli random vector. To upper bound the tail probability on the right-hand side, we resort to Hoeffding's inequality, which gives

$$\mathbb{P} \left( \sum_{i=1}^n \epsilon_i - \frac{n}{2} > \frac{n}{2} \left(1 - \left(\frac{1}{2}\right)^{p-1}\right) \right) \leq \exp \left( -\frac{n}{2} \left(1 - \left(\frac{1}{2}\right)^{p-1}\right)^2 \right).$$

By transitivity, this implies that a sufficient condition for

$$\mathbb{P} \left( \exists (j, j') \in \llbracket 1; q \rrbracket^2 : 1 \leq j < j' \leq q \text{ and } \sum_{i=1}^n \mathbf{1}_{\{\epsilon_{j,i} \neq \epsilon_{j',i}\}} < \left(\frac{1}{2}\right)^p n \right) < 1$$

is

$$\binom{q}{2} \exp \left( -\frac{n}{2} \left(1 - \left(\frac{1}{2}\right)^{p-1}\right)^2 \right) < 1$$

and consequently

$$q \leq \left\lceil 2^{\frac{n}{K^p}} \right\rceil,$$

which is precisely the value announced and thus concludes the proof.  $\square$

The transition between covering and packing numbers is provided by a well-known equivalence.

**Lemma 22.** *Let  $(\mathcal{E}, \rho)$  be a pseudo-metric space. For every totally bounded set  $\mathcal{E}' \subset \mathcal{E}$  and  $\epsilon \in \mathbb{R}_+^*$ ,  $\mathcal{M}(2\epsilon, \mathcal{E}', \rho) \leq \mathcal{N}^{int}(\epsilon, \mathcal{E}', \rho) \leq \mathcal{M}(\epsilon, \mathcal{E}', \rho)$ .*

## C.2 Margin Natarajan Dimension of $\rho_G$

The proof of Lemma 10 is the following one.

*Proof.* Suppose that for  $\gamma \in (0, M_G]$ , the triplet  $(s_{Z^n}, \mathbf{b}_n, \mathbf{c}_n)$  is  $\gamma$ - $N$ -shattered by  $\rho_G$ . According to Proposition 2, in order to upper bound  $n$ , i.e.,  $\gamma$ - $N$ -dim  $(\rho_G)$ , one can assume that for every  $i \in \llbracket 1; n \rrbracket$ ,  $y_i < c_i$ , and the biases can be negative. Let  $\bar{\mathcal{G}} = \{g^{\mathbf{s}_n} : \mathbf{s}_n \in \{-1, 1\}^n\}$

be a subset of  $\mathcal{G}$  (of minimal cardinality) such that  $\rho_{\bar{\mathcal{G}}} = \{\rho_{g^{\mathbf{s}_n}} : \mathbf{s}_n \in \{-1, 1\}^n\}$   $\gamma$ - $N$ -shatters  $(s_{\mathcal{Z}^n}, \mathbf{b}_n, \mathbf{c}_n)$ . For every pair  $(k, l) \in \llbracket 1; C \rrbracket^2$  satisfying  $k < l$ , let  $\mathcal{S}_{k,l}$  be the subset of  $\llbracket 1; n \rrbracket$  defined as follows:

$$\mathcal{S}_{k,l} = \{i \in \llbracket 1; n \rrbracket : y_i = k \text{ and } c_i = l\}$$

and let  $n_{k,l} \in \llbracket 0; n \rrbracket$  be its cardinality. By construction,  $\mathcal{P} = \{\mathcal{S}_{k,l} : n_{k,l} > 0\}$  is a partition of  $\llbracket 1; n \rrbracket$ . For every vector  $\mathbf{s}_n = (s_i)_{1 \leq i \leq n} \in \{-1, 1\}^n$ , the function  $g^{\mathbf{s}_n}$  satisfies:

$$\forall i \in \llbracket 1; n \rrbracket, \begin{cases} \text{if } s_i = 1, \rho_{g^{\mathbf{s}_n}}(x_i, y_i) - b_i \geq \gamma \\ \text{if } s_i = -1, \rho_{g^{\mathbf{s}_n}}(x_i, c_i) + b_i \geq \gamma \end{cases}.$$

For a fixed  $\mathcal{S}_{k,l} \in \mathcal{P}$ , this implies that

$$\forall i \in \mathcal{S}_{k,l}, s_i \left( \frac{1}{2} (g_k^{\mathbf{s}_n}(x_i) - g_l^{\mathbf{s}_n}(x_i)) - b_i \right) \geq \gamma.$$

Let  $\mathcal{D}_{\mathcal{G},k,l} = \{\frac{1}{2}(g_k - g_l) : g \in \mathcal{G}\}$ . By definition, we have established that its subset  $\{\frac{1}{2}(g_k^{\mathbf{s}_n} - g_l^{\mathbf{s}_n}) : \mathbf{s}_n \in \{-1, 1\}^n\}$   $\gamma$ -shatters a set of cardinality  $n_{k,l}$ , with the consequence that

$$n_{k,l} \leq \gamma\text{-dim}(\mathcal{D}_{\mathcal{G},k,l}).$$

Summing over all the elements of the partition  $\mathcal{P}$  gives

$$\gamma\text{-N-dim}(\rho_{\mathcal{G}}) \leq \sum_{k=1}^{C-1} \sum_{l=k+1}^C \gamma\text{-dim}(\mathcal{D}_{\mathcal{G},k,l}). \quad (75)$$

Inequality (75) implies (15) since by definition,  $\mathcal{D}_{\mathcal{G}} = \bigcup_{1 \leq k < l \leq C} \mathcal{D}_{\mathcal{G},k,l}$ . The second upper bound on the  $\gamma$ -dimensions of the classes  $\mathcal{D}_{\mathcal{G},k,l}$  is obtained by application of the standard strategy outlined in Section 3.1. Applying in sequence Proposition 3 and Lemma 22 (left-hand side inequality) gives:

$$\begin{aligned} \gamma\text{-dim}(\mathcal{D}_{\mathcal{G},k,l}) &\leq 16 \log_2(\mathcal{M}_2(\gamma, \mathcal{D}_{\mathcal{G},k,l}, \gamma\text{-dim}(\mathcal{D}_{\mathcal{G},k,l}))) \\ &\leq 16 \log_2\left(\mathcal{N}_2^{\text{int}}\left(\frac{\gamma}{2}, \mathcal{D}_{\mathcal{G},k,l}, \gamma\text{-dim}(\mathcal{D}_{\mathcal{G},k,l})\right)\right). \end{aligned}$$

An optimization of the proof of Lemma 2 in the degenerate case  $C = 2$  produces

$$\mathcal{N}_2^{\text{int}}\left(\frac{\gamma}{2}, \mathcal{D}_{\mathcal{G},k,l}, \gamma\text{-dim}(\mathcal{D}_{\mathcal{G},k,l})\right) \leq \mathcal{N}_2^{\text{int}}\left(\frac{\gamma}{2}, \mathcal{G}_k, \gamma\text{-dim}(\mathcal{D}_{\mathcal{G},k,l})\right) \times \mathcal{N}_2^{\text{int}}\left(\frac{\gamma}{2}, \mathcal{G}_l, \gamma\text{-dim}(\mathcal{D}_{\mathcal{G},k,l})\right).$$

To upper bound the covering numbers in the right-hand side, it suffices to apply Lemma 22 (right-hand side inequality) and Theorem 1 in Mendelson and Vershynin (2003) (with the

optimized constants of Lemma 7). This produces:

$$\begin{aligned} \mathcal{N}_2^{\text{int}} \left( \frac{\gamma}{2}, \mathcal{G}_k, \gamma\text{-dim}(\mathcal{D}_{\mathcal{G},k,l}) \right) &\leq \mathcal{M}_2 \left( \frac{\gamma}{2}, \mathcal{G}_k, \gamma\text{-dim}(\mathcal{D}_{\mathcal{G},k,l}) \right) \\ &\leq \left( \frac{20M_{\mathcal{G}}}{\gamma} \right)^{12d_k(\frac{\gamma}{48})}, \end{aligned}$$

where  $d_k(\epsilon) = \epsilon\text{-dim}(\mathcal{G}_k)$ . The second upper bound on the  $\gamma$ -dimensions of the classes  $\mathcal{D}_{\mathcal{G},k,l}$  has been obtained. By substitution into (75), we get

$$\begin{aligned} \gamma\text{-N-dim}(\rho_{\mathcal{G}}) &\leq 192(C-1) \log_2 \left( \frac{20M_{\mathcal{G}}}{\gamma} \right) \sum_{k=1}^C \left( \frac{\gamma}{48} \right)\text{-dim}(\mathcal{G}_k) \\ &\leq 384 \binom{C}{2} \log_2 \left( \frac{20M_{\mathcal{G}}}{\gamma} \right) \left( \frac{\gamma}{48} \right)\text{-dim}(\mathcal{G}_0). \end{aligned}$$

□

### C.3 Margin Natarajan Dimension of $\rho_{\mathcal{H}_{\Lambda}}$

The proof of Lemma 11 reuses the notations of the proof of Lemma 10, with  $\mathcal{G}$  set equal to  $\mathcal{H}_{\Lambda}$ .

*Proof.* By application of Khintchine inequality, every set  $\mathcal{S}_{k,l}$  can be split into two subsets  $\mathcal{S}_{k,l}^+$  and  $\mathcal{S}_{k,l}^-$  such that

$$\left\| \sum_{i \in \mathcal{S}_{k,l}^+} \kappa_{x_i} - \sum_{i \in \mathcal{S}_{k,l}^-} \kappa_{x_i} \right\|_{\mathbf{H}_{\kappa}} \leq \sqrt{n_{k,l}} \Lambda_{\mathcal{X}}. \quad (76)$$

Let the vector  $\mathbf{s}_n \in \{-1, 1\}^n$  be defined as follows:

$$\forall i \in \llbracket 1; n \rrbracket, \begin{cases} i \in \mathcal{S}_{k,l}^+ \implies s_i = 1 \\ i \in \mathcal{S}_{k,l}^- \implies s_i = -1 \end{cases}.$$

For every  $\mathcal{S}_{k,l} \in \mathcal{P}$  and  $i \in \mathcal{S}_{k,l}$ , applying the reproducing property gives

$$\begin{cases} \text{if } s_i = 1, \frac{1}{2} \langle h_k^{\mathbf{s}_n} - h_l^{\mathbf{s}_n}, \kappa_{x_i} \rangle_{\mathbf{H}_{\kappa}} - b_i \geq \gamma \text{ and } \frac{1}{2} \langle h_l^{-\mathbf{s}_n} - h_k^{-\mathbf{s}_n}, \kappa_{x_i} \rangle_{\mathbf{H}_{\kappa}} + b_i \geq \gamma \\ \text{if } s_i = -1, \frac{1}{2} \langle h_l^{\mathbf{s}_n} - h_k^{\mathbf{s}_n}, \kappa_{x_i} \rangle_{\mathbf{H}_{\kappa}} + b_i \geq \gamma \text{ and } \frac{1}{2} \langle h_k^{-\mathbf{s}_n} - h_l^{-\mathbf{s}_n}, \kappa_{x_i} \rangle_{\mathbf{H}_{\kappa}} - b_i \geq \gamma \end{cases}. \quad (77)$$

For every  $k \in \llbracket 1; C \rrbracket$ , let  $h_k^{\delta} = \frac{1}{2} (h_k^{\mathbf{s}_n} - h_k^{-\mathbf{s}_n})$ . Then (77) produces by summation:

$$\forall i \in \mathcal{S}_{k,l}, \begin{cases} i \in \mathcal{S}_{k,l}^+ \implies \frac{1}{2} \langle h_k^{\delta} - h_l^{\delta}, \kappa_{x_i} \rangle_{\mathbf{H}_{\kappa}} \geq \gamma \\ i \in \mathcal{S}_{k,l}^- \implies \frac{1}{2} \langle h_k^{\delta} - h_l^{\delta}, -\kappa_{x_i} \rangle_{\mathbf{H}_{\kappa}} \geq \gamma \end{cases}. \quad (78)$$

By summation over  $i \in \mathcal{S}_{k,l}$ , it results from (78) that:

$$\frac{1}{2} \left\langle h_k^\delta - h_l^\delta, \sum_{i \in \mathcal{S}_{k,l}^+} \kappa_{x_i} - \sum_{i \in \mathcal{S}_{k,l}^-} \kappa_{x_i} \right\rangle_{\mathbf{H}_\kappa} \geq n_{k,l} \gamma. \quad (79)$$

Applying the Cauchy-Schwarz inequality to (79) yields

$$\frac{1}{2} \left\| h_k^\delta - h_l^\delta \right\|_{\mathbf{H}_\kappa} \left\| \sum_{i \in \mathcal{S}_{k,l}^+} \kappa_{x_i} - \sum_{i \in \mathcal{S}_{k,l}^-} \kappa_{x_i} \right\|_{\mathbf{H}_\kappa} \geq n_{k,l} \gamma. \quad (80)$$

By substitution of (76) into (80),

$$\forall \mathcal{S}_{k,l} \in \mathcal{P}, \quad n_{k,l} \leq \left( \frac{\frac{1}{2} \left\| h_k^\delta - h_l^\delta \right\|_{\mathbf{H}_\kappa} \Lambda \mathcal{X}}{\gamma} \right)^2.$$

By summation over all the elements of the partition  $\mathcal{P}$ ,

$$n \leq \left( \frac{\Lambda \mathcal{X}}{2\gamma} \right)^2 \sum_{1 \leq k < l \leq C} \left\| h_k^\delta - h_l^\delta \right\|_{\mathbf{H}_\kappa}^2. \quad (81)$$

To upper bound the sum in the right-hand side of (81), we first note that

$$\left\| h_k^\delta - h_l^\delta \right\|_{\mathbf{H}_\kappa}^2 \leq \frac{1}{2} \left( \left\| h_k^{s_n} - h_l^{s_n} \right\|_{\mathbf{H}_\kappa}^2 + \left\| h_k^{-s_n} - h_l^{-s_n} \right\|_{\mathbf{H}_\kappa}^2 \right),$$

with the consequence that

$$n \leq \left( \frac{\Lambda \mathcal{X}}{2\gamma} \right)^2 \frac{1}{2} \sum_{1 \leq k < l \leq C} \left( \left\| h_k^{s_n} - h_l^{s_n} \right\|_{\mathbf{H}_\kappa}^2 + \left\| h_k^{-s_n} - h_l^{-s_n} \right\|_{\mathbf{H}_\kappa}^2 \right). \quad (82)$$

Now, since by hypothesis,  $\sum_{k=1}^C h_k = \mathbf{0}_{\mathbf{H}_\kappa}$ ,

$$\begin{aligned} \sum_{1 \leq k < l \leq C} \left\| h_k^{s_n} - h_l^{s_n} \right\|_{\mathbf{H}_\kappa}^2 &= (C-1) \sum_{k=1}^C \left\| h_k^{s_n} \right\|_{\mathbf{H}_\kappa}^2 - 2 \sum_{1 \leq k < l \leq C} \langle h_k^{s_n}, h_l^{s_n} \rangle_{\mathbf{H}_\kappa} \\ &= C \sum_{k=1}^C \left\| h_k^{s_n} \right\|_{\mathbf{H}_\kappa}^2 - \sum_{k=1}^C \sum_{l=1}^C \langle h_k^{s_n}, h_l^{s_n} \rangle_{\mathbf{H}_\kappa} \\ &= C \sum_{k=1}^C \left\| h_k^{s_n} \right\|_{\mathbf{H}_\kappa}^2 - \left\| \sum_{k=1}^C h_k^{s_n} \right\|_{\mathbf{H}_\kappa}^2 \\ &= C \sum_{k=1}^C \left\| h_k^{s_n} \right\|_{\mathbf{H}_\kappa}^2 \\ &= C \left\| h^{s_n} \right\|_{\mathbf{H}_{\kappa,C}}^2 \\ &\leq C \Lambda^2. \end{aligned} \quad (83)$$

Obviously, (83) also provides an upper bound on  $\sum_{1 \leq k < l \leq C} \left\| h_k^{-s_n} - h_l^{-s_n} \right\|_{\mathbf{H}_\kappa}^2$ . Thus, a substitution into (82) concludes the proof.  $\square$

## D Upper Bound on the Rademacher Complexity

The proof of Theorem 3 is the following one.

*Proof.* In all three cases, the starting point is Inequality (28).

**First case:**  $d_{\mathcal{G},\gamma} \in (0, 2)$

This case is the only one for which the entropy integral exists. Setting for every  $j \in \mathbb{N}$ ,  $h(j) = \gamma 2^{-2-d_{\mathcal{G},\gamma}j}$ , we obtain

$$R_m(\rho_{\mathcal{G},\gamma}) \leq 24 \left(1 + 2^{\frac{2}{2-d_{\mathcal{G},\gamma}}}\right) \sqrt{\frac{F_1(C)}{m}} \gamma^{1-\frac{d_{\mathcal{G},\gamma}}{2}} \int_0^{\frac{1}{2}} \ln \left( (C-1) \left(4\epsilon^{-\frac{2}{2-d_{\mathcal{G},\gamma}}}\right)^5 \right) d\epsilon.$$

Let  $I(C)$  denote the integral. Then,

$$\begin{aligned} I(C) &= \int_0^{\frac{1}{2}} \ln \left( (C-1) \left(4\epsilon^{-\frac{2}{2-d_{\mathcal{G},\gamma}}}\right)^5 \right) d\epsilon \\ &= \frac{1}{2} \left( \ln((C-1)4^5) + 10 \frac{1 + \ln(2)}{2 - d_{\mathcal{G},\gamma}} \right). \end{aligned}$$

**Second case:**  $d_{\mathcal{G},\gamma} = 2$

$$R_m(\rho_{\mathcal{G},\gamma}) \leq h(N) + \frac{8}{3} \sqrt{\frac{F_1(C)}{m}} \sum_{j \in \mathcal{J}} \frac{h(j) + h(j-1)}{h(j)} \log_2 \left( (C-1) \left( \frac{4\gamma}{h(j)} \right)^5 \right).$$

We set  $N = \left\lceil \log_2 \left( \frac{\sqrt{m}}{\log_2(m)} \right) \right\rceil$  and  $h(j) = \gamma \frac{\log_2(m)}{\sqrt{m}} 2^{N-j}$ . Note that  $m \geq 2$  ensures that  $N \geq 1$ . Then,

$$\begin{aligned} R_m(\rho_{\mathcal{G},\gamma}) &\leq \gamma \frac{\log_2(m)}{\sqrt{m}} + 8 \sqrt{\frac{F_1(C)}{m}} \sum_{j=1}^N \log_2 \left( (C-1) \left( \frac{4\gamma}{h(j)} \right)^5 \right) \\ &\leq \gamma \frac{\log_2(m)}{\sqrt{m}} + 8 \sqrt{\frac{F_1(C)}{m}} \left\lceil \log_2 \left( \frac{\sqrt{m}}{\log_2(m)} \right) \right\rceil \log_2 \left( (C-1) \left( 4 \frac{\sqrt{m}}{\log_2(m)} \right)^5 \right). \end{aligned}$$

**Third case:**  $d_{\mathcal{G},\gamma} > 2$

For  $N = \left\lceil \frac{d_{\mathcal{G},\gamma}-2}{2d_{\mathcal{G},\gamma}} \log_2 \left( \frac{m}{\log_2(m)} \right) \right\rceil$ , let us set  $h(j) = \gamma \left( \frac{\log_2(m)}{m} \right)^{\frac{1}{d_{\mathcal{G},\gamma}}} 2^{\frac{2}{d_{\mathcal{G},\gamma}-2}(-j+N)}$ . We then get

$$R_m(\rho_{\mathcal{G},\gamma}) \leq \gamma \left( \frac{\log_2(m)}{m} \right)^{\frac{1}{d_{\mathcal{G},\gamma}}} \left[ 1 + \frac{8}{3} \left( 1 + 2^{\frac{2}{d_{\mathcal{G},\gamma}-2}} \right) \left( \frac{1}{\gamma} \right)^{\frac{d_{\mathcal{G},\gamma}}{2}} \sqrt{\frac{F_1(C)}{\log_2(m)}} S_N \right]$$

with

$$\begin{aligned} S_N &= \sum_{j=1}^N 2^{j-N} \log_2 \left( (C-1) \left( \frac{4\gamma}{h(j)} \right)^5 \right) \\ &\leq \log_2 \left( (C-1) \left( \frac{4\gamma}{h(N)} \right)^5 \right) \sum_{j=1}^N 2^{j-N} \\ &\leq 2 \log_2 \left( (C-1) \left( 4 \left( \frac{m}{\log_2(m)} \right)^{\frac{1}{d_{\mathcal{G},\gamma}}} \right)^5 \right). \end{aligned}$$

□