# MINIMAX LEARNING RATES FOR ESTIMATING BINARY CLASSIFIERS UNDER MARGIN CONDITIONS

## **Anonymous authors**

Paper under double-blind review

#### **ABSTRACT**

We study classification problems using binary estimators where the decision boundary is described by horizon functions and where the data distribution satisfies a geometric margin condition. We establish lower bounds for the minimax learning rate over broad function classes with bounded Kolmogorov entropy in Lebesgue norms. A key novelty of our work is the derivation of lower bounds on the worst-case learning rates under a geometric margin condition—a setting that is almost universally satisfied in practice but remains theoretically challenging. Moreover, our results deal with the noiseless setting, where lower bounds are particularly hard to establish. We apply our general results to classification problems with decision boundaries belonging to several function classes: for Barron-regular functions, Hölder-continuous functions, and convex functions with strong margins, we identify optimal rates close to the fast learning rates of  $\mathcal{O}(n^{-1})$  for  $n \in \mathbb{N}$  samples.

#### 1 Introduction

How well can we solve classification problems with complex decision boundaries in deep learning? A lot of emphasis has been put on the noise in the problem. However, in practice, data sets may have very strong margins (see C3)) between the classes, which makes learning much simpler (see e.g. Figure 1). The presence of a margin seriously complicates the identification of lower bounds on learning success. This is intuitively clear, since in the extreme case, where certain regions between the classes almost surely do not contain any data points, many decision boundaries are valid.

In this manuscript, we overcome these issues and present lower bounds on learning under margin conditions.

For  $n \in \mathbb{N}$  samples, an estimator can be defined as a measurable function  $f: \Lambda^n \to \mathcal{F}$  and a binary classifier can be seen as an indicator function  $\mathbb{1}_{\Omega}: X \to \{0,1\}$ , where  $\Lambda:=X \times \{0,1\}$  for some set X containing a sequence  $\{x_i\}_{i=1}^n$  that belongs to a sample  $((x_i,\mathbb{1}_{\Omega}(x_i)))_{i=1}^n \in \Lambda^n$ ,  $\mathcal{F}$  is a measurable space, and  $\Omega \subset X$  is the decision set.

**Remark 1.** It is clear from the definition above, that the classifiers considered in this work are not corrupted by noise. This is an important assumption if we want to resolve the precise role of the regularity of the decision boundary and the margin conditions. Indeed, even the presence of low noise could yield vastly different lower bounds, because the learning problem then requires resolving the noise and this complication can mask the role of the decision boundary and the margin. An extended discussion of this is given in (Petersen & Voigtlaender, 2021, Section 1.1, Point 1). For a quick argument, we highlight, e.g., Stone (1982), where it was obtained that the optimal learning rate, to learn a function  $f \in C^k([0,1]^d)$  with  $\|f\|_{C^k} \leq 1$  and noise defined as a parameter  $\varepsilon \stackrel{iid}{\sim} N(0,\sigma^2)$  for  $\sigma > 0$ , is of the order of  $O(n^{-k/(2k+d)})$ , and decays slower than  $n^{-1/2}$ . On the other hand, in Krieg & Sonnleitner (2023), where the same problem is considered without noise, learning rates of the order of  $O(n^{-k/d})$  were obtained, in some cases faster than  $n^{-1}$ .

## 1.1 CONDITIONS

In this paper, we study classification learning problems using binary estimators, when  $\mathcal{F} := L^2([0,1]^d)$ ;  $X := [0,1]^d$  and the following conditions are fulfilled:

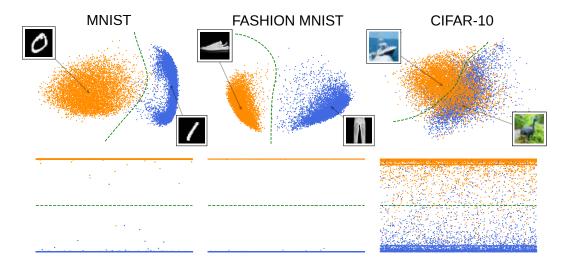


Figure 1: Geometric margin in common classification problems. The top row shows a two dimensional embedding on the first two principle components and a decision boundary identified by a support vector machine. Clearly MNIST Lecun et al. (1998) and Fashion MNIST Xiao et al. (2017) exhibit a strong margin between some classes. For the CIFAR-10 data Krizhevsky (2009) the margin is not visible in the two dimensional embedding. In the second row, we show the class probabilities predicted by a support vector classifier, which again shows extremely strong margin for MNIST and Fashon MNIST, but also reveals that the CIFAR-10 data exhibits a margin, albeit a weaker one. Which lower bounds on learning can be found in the presence of such various types of margins will be demonstrated in our main results Theorem 6, Corollary 8 and Corollary 9.

C1)  $\Omega$  can be described by horizon functions: We define the general horizon function associated with  $b \in \mathscr{C} \subset C([0,1]^{d-1};[0,1])$  as

$$h_b: [0,1]^d \to \{0,1\},$$
  
 $\mathbf{x} = (x_1, \dots, x_d) \mapsto \mathbb{1}_{b(\mathbf{x}^{(d)}) \le x_d},$  (1)

where  $\boldsymbol{x}^{(d)} := (x_1, \dots, x_{d-1})$ . Then it is fulfilled that

$$\Omega = \Omega(h) := \Omega_h = \left\{ \boldsymbol{x} \in [0, 1]^d : h(\boldsymbol{x}) = 1 \right\} \quad \text{with} \quad h \in H_{\mathscr{C}} := \left\{ h_b : b \in \mathscr{C} \right\}, \quad (2)$$

where  $H_{\mathscr{C}}$  is the set of general horizon functions associated to  $\mathscr{C}$  Petersen & Voigtlaender (2021). So,  $\mathbb{1}_{\Omega_h} = h$ .

Moreover,  $\Lambda$  is equipped with a probability measure  $\mu_h$ , such that  $\mu_h\left([0,1]^d\times\{i\}\right)=\nu_i\in[0,1]$ , for  $i\in\{0,1\}$  and  $\nu_0+\nu_1=1$ ;  $\{\boldsymbol{x}_i\}_{i=1}^n\stackrel{iid}{\sim}\mu_h$ , where  $h\in H_{\mathscr{C}}$ , and  $\mu_h$  is the marginal distribution of  $[0,1]^d$  admitting a density function  $f_h$  with respect to the Lebesgue measure  $\lambda$ .

C2) Regular boundary:  $\mathscr C$  is convex with  $0 \in \mathscr C$ . Furthermore, every  $b \in \mathscr C$  satisfies

$$|b(\boldsymbol{z}) - b(\boldsymbol{z}')| \le K \|\boldsymbol{z} - \boldsymbol{z}'\|^{\alpha}$$
 for all  $\boldsymbol{z}, \boldsymbol{z}' \in [0, 1]^{d-1}$ ,

some  $\alpha \in (0,1]$  and a constant K > 0.

C3) The margin condition is satisfied: There exist  $C, \gamma > 0$ , such that for all  $\epsilon > 0$ ,

$$\mu_h\left(B_{\epsilon}^h\right) \leq C\epsilon^{\gamma}$$
 where  $B_{\epsilon}^h := \left\{ \boldsymbol{x} \in [0,1]^d : \operatorname{dist}\left(\boldsymbol{x}, \partial \Omega_h\right) \leq \epsilon \right\}$ 

is the ball of radius  $\epsilon>0$  around  $\partial\Omega_h$  with respect to the Euclidean distance

$$\operatorname{dist}(\boldsymbol{x},\partial\Omega_h) := \inf_{\boldsymbol{x'} \in \partial\Omega_h} \|\boldsymbol{x} - \boldsymbol{x'}\|,$$

and  $\gamma$  is called the margin exponent Christmann & Steinwart (2008); Kim et al. (2021).

Additionally, we define the measure  $\mu$  on  $\Lambda$  with marginal  $\lambda$  on  $[0,1]^d$  and

$$\mu([0,1]^d \times \{i\}) = \mu_h([0,1]^d \times \{i\}) = \nu_i \text{ for all } i \in \{0,1\}.$$
 (3)

Under the above conditions, we establish lower bounds for the minimax error associated with the problem of estimating binary classifiers when the learning set satisfies C1) on a function class  $\mathscr C$  with regularity C2), and the data distribution satisfies the margin condition C3); i.e., we lower bound the following inf-sup expression

$$\mathcal{I}_{n}(\mathscr{C}) := \inf_{A \in \mathcal{A}_{n}(L^{2}(\lambda))} \sup_{\substack{h \in H_{\mathscr{C}} \\ \mu_{h} \text{ satisfies C3)}}} \mathbb{E}_{\{\boldsymbol{x}_{i}\}_{i=1}^{n} \stackrel{iid}{\sim} \mu_{h}} \|A(\boldsymbol{S}_{h}) - h\|_{L^{2}(\mu_{h})}^{2}, \tag{4}$$

where  $S_h := ((x_i, h(x_i)))_{i=1}^n$  is the sample of size  $n \in \mathbb{N}$ , for a function  $h : [0, 1]^d \to \{0, 1\}$ , and

$$\mathcal{A}_n(\mathcal{G}) := \{A : \Lambda^n \to \mathcal{G} : A \text{ is measurable}\} \quad \text{ with } \quad \mathcal{G} \subseteq L^2(\lambda) := L^2([0,1]^d, \lambda).$$

Remark 2. In particular, the margin condition C3) is satisfied when

$$f_h(\boldsymbol{x}) \lesssim \begin{cases} \min\left\{\frac{\epsilon^{\gamma}}{\lambda(B_{\epsilon}^h)}, 1\right\} & \textit{if } \boldsymbol{x} \in B_{\epsilon}^h \\ 1 & \textit{otherwise,} \end{cases}$$

for almost every  $x \in [0,1]^d$  and for all  $\epsilon > 0$ . Moreover, C3) holds when

$$f_h(\boldsymbol{x}) \lesssim \operatorname{dist}^{\gamma}(\boldsymbol{x}, \partial \Omega_h)$$
 for almost every  $[0, 1]^d$ .

#### 1.2 Entropy

By Hurewicz & Wallman (1948); Kolmogorov & Tihomirov (1993), it is known that it is possible to characterize the "mass" of sets in metric spaces through the scaling law of the number of balls necessary to cover them with variable radii. Moreover, this characterization can be extended to spaces of densities with respect to distances that are not necessarily a metric, such as the distance associated with the Kullback-Leibler divergence or the Hellinger distance (see Yang & Barron (1999)). This leads to the definition of the concept of  $\varepsilon$ -covering entropy (or  $\varepsilon$ -entropy of Kolmogorov) which serves as a measure of complexity for sets, where a complex set is one that needs many elements to cover it and a simple set does not.

Let  $\mathcal{X} \neq \emptyset$  be a set and let  $\rho: \mathcal{X} \times \mathcal{X} \to [0, \infty]$  be a distance function (without the need to be symmetric or satisfy some form of triangle inequality). For a set  $\emptyset \neq K \subset \mathcal{X}$  and  $\varepsilon > 0$ , it is said that:

(N)  $G(\varepsilon) \subseteq \mathcal{X}$  is an  $\varepsilon$ -net of K, when for all  $x \in K$ , there exists  $y \in G(\varepsilon) \subset \mathcal{X}$ , such that  $\rho(x,y) \leq \varepsilon$ .

We define

$$V_{K,\rho}(\varepsilon) := \log \min \{ |G_{\varepsilon}| : G_{\varepsilon} \text{ satisfies (N)} \}$$
 as the  $\varepsilon$ -covering entropy of  $K$ . (5)

Moreover, a similar notion can be defined as follows.

(P)  $G(\varepsilon) \subseteq K$  is an  $\varepsilon$ -packing (or  $\varepsilon$ -distinguishable) set in K, when for all  $x_1, x_2 \in G(\varepsilon)$ , it holds that  $\rho(x_1, x_2) > \varepsilon$ .

We define (the  $\varepsilon$ -capacity in Kolmogorov & Tihomirov (1993))

$$M_{K,o}(\varepsilon) := \log \max \{ |G_{\varepsilon}| : G_{\varepsilon} \text{ satisfies (P)} \}$$
 as the  $\varepsilon$ -packing entropy of  $K$ . (6)

Indeed, when  $\rho$  is a metric, (5) and (6) are equivalent in the sense that (see Kolmogorov & Tihomirov (1993) and (Petersen & Voigtlaender, 2021, Remark 3.10))

$$M_{K,\rho}(2\varepsilon) \le V_{K,\rho}(\varepsilon) \le M_{K,\rho}(\varepsilon).$$
 (7)

To conclude, we introduce the following distance functions between probability densities that will be applied to the notions of entropy. We define the Hellinger distance  $\rho_H$  and the associated distance of the Kullback-Leibler divergence  $\rho_{KL}$ , as follows. Let

$$D_{m{\mu}} := \left\{ f : \Lambda o [0,\infty) : f ext{ measurable and } \int_{\Lambda} f \, dm{\mu} = 1 
ight\}$$

be the set of all probability densities functions (with respect to  $\mu$ ) on  $\Lambda$ . The distances  $\rho_H$  and  $\rho_{KL}$ , between  $p, q \in D_{\mu}$  are defined as

$$\rho_H(p,q) := \left( \int_{\Lambda} (\sqrt{p} - \sqrt{q})^2 d\boldsymbol{\mu} \right)^{1/2} \quad \text{and} \quad \rho_{KL}(p,q) := \left( \int_{\Lambda} p \log(p/q) d\boldsymbol{\mu} \right)^{1/2}. \tag{8}$$

#### 1.3 Entropy bounds for some function classes

We introduce some results on entropy bounds under the  $L^p$  norm with  $p \ge 1$ , for classes of functions previously studied in Guntuboyina & Sen (2013), Kolmogorov & Tihomirov (1993) and Petersen & Voigtlaender (2021). Our goal is to apply Theorem 6 to classes such as  $\mathscr{C} \subset C([0,1]^{d-1};[0,1])$  in C1).

# 1.3.1 BARRON REGULAR FUNCTIONS

A function  $f:[0,1]^{d-1}\to\mathbb{R}$  is said to be of Barron class (see Barron (1993); Caragea et al. (2023); Petersen & Voigtlaender (2021)) with constant C>0, if there are  $c\in[-C,C]$  and a measurable function  $F:\mathbb{R}^{d-1}\to\mathbb{C}$  satisfying

$$\int_{\mathbb{R}^{d-1}} |F(\boldsymbol{\xi})| \sup_{\boldsymbol{x} \in [0,1]^{d-1}} |\langle \boldsymbol{\xi}, \boldsymbol{x} \rangle| \, d\boldsymbol{\xi} \le C \quad \text{and} \quad f(\boldsymbol{x}) = c + \int_{\mathbb{R}^{d-1}} (e^{i\langle \boldsymbol{x}, \boldsymbol{\xi} \rangle} - 1) \cdot F(\boldsymbol{\xi}) d\boldsymbol{\xi}$$
(9)

for all  $x \in [0,1]^{d-1}$ . The set of all these Barron functions is denoted as  $\mathcal{B}_C$ . Propositions 4.4 and 4.6 in Petersen & Voigtlaender (2021) state the following lemma.

**Lemma 3.** For all  $0 < \varepsilon < 1$ ,

$$\varepsilon^{-\frac{2(d-1)}{d+1}} \lesssim M_{\mathcal{B}_C, L^1([0,1]^{d-1})}(\varepsilon) \leq M_{\mathcal{B}_C, L^{\infty}([0,1]^{d-1})}(\varepsilon) \lesssim \varepsilon^{-\frac{2(d-1)}{d+1}}(1 + \log(1/\varepsilon)).$$

## 1.3.2 $\alpha$ -Hölder continuous functions

According to Clements (1963) and (Kolmogorov & Tihomirov, 1993, Sections  $\S 5, \S 9$ ), let  $\mathcal{H}_{\alpha}$  be the class of functions f from  $[0,1]^{d-1}$  to [0,1], that are bounded by a constant C and have all partial derivatives of the order  $i \le k := \beta - \alpha \in \mathbb{Z}_{\ge 0}$  with the kth order derivatives satisfying the  $\alpha$ -Hölder condition with exponent  $\alpha \in (0,1]$ , i.e.

$$\left| f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2) \right| \le L \left\| \boldsymbol{x}_1 - \boldsymbol{x}_2 \right\|_{\infty}^{\alpha} := L \left[ \max_{j \in \{1, ..., d-1\}} \left| x_{1,j} - x_{2,j} \right| \right]^{\alpha},$$

for some constant L>0, where  $\boldsymbol{x}_i:=(x_{i,1},\ldots,x_{i,d-1})\in[0,1]^{d-1}$ , for  $i\in\{1,2\}$ . Then, Clements (1963) and (Kolmogorov & Tihomirov, 1993,  $\S 9$ -1) and Theorem V) imply the next result.

**Lemma 4.** For all 
$$\alpha > 0$$
 and  $p \in \{1, \infty\}$ ,  $V_{\mathcal{H}_{\alpha}, L^{p}([0,1]^{d-1})}(\varepsilon) \approx \varepsilon^{-(d-1)/\alpha}$ .

## 1.3.3 Convex functions

We denote by  $\mathcal{C}_B([a,b]^{d-1})$  the set of all convex functions on  $[a,b]^{d-1}$  with a < b and  $d \geq 2$ , that are uniformly bounded by B. For this class of functions, the  $\varepsilon$ -covering entropy has been studied when d=2 in Dryanov (2009); by adding an uniformly Lipschitz condition in Bronshtein (1976); and without extra conditions for all  $d \geq 2$ , in Guntuboyina & Sen (2013). We present the following lemma, which is a consequence of Theorems 3.1 and 3.3 in Guntuboyina & Sen (2013), for the case that a=0 and b=1.

**Lemma 5.** For all  $p \in [1, \infty)$ , B > 0 and  $0 < \varepsilon \lesssim B$ , it is fulfilled that

$$V_{\mathcal{C}_B,L^p([0,1]^{d-1})}(\varepsilon) \approx \varepsilon^{-(d-1)/2},$$

where  $C_B := C_B([0,1]^{d-1})$  and  $d \ge 2$ .

#### 1.4 Previous works and our contribution

Some related work on binary classifiers under the margin condition C3) can be found in: (Christmann & Steinwart, 2008, Section 8), for SVMs where learning rates were sometimes as fast as  $n^{-1}$ , being n the number of data points; Kim et al. (2021), based on neural networks with hinge loss that achieves fast convergence rates when  $d \lesssim \gamma$ , particularly as fast as  $n^{-(q+1)/(q+2)}$  when the margin exponent  $\gamma \to \infty$ , where q is a noise parameter; and García & Petersen (2025), where it is found that using ReLU neural networks, the strong margin conditions imply fast learning bounds that are close to  $n^{-1}(1 + \log n)$ . Furthermore, in García & Petersen (2025); Kim et al. (2021), a regularity condition is assumed on the decision boundary, where  $\partial\Omega$  can be described by classes of functions  $\mathscr{C}$  satisfying condition C1) on the elements of some covering for the set  $\Omega$ . In Kim et al. (2021), the class of functions  $\mathscr{C} \subset \mathcal{H}_{\alpha}$  is considered, and in García & Petersen (2025), the class  $\mathscr{C} = \mathcal{B}_{C}$ . The last two works mentioned above, only provide upper bounds for the learning rate when binary classifiers are approximated by neural networks under the margin condition. Our contribution now, by finding minimax lower bounds for learning rates on binary estimators, is to confirm that these learning rates are indeed optimal. It is important to highlight that the margin plays a main role in the fast learning rates obtained in each function space, when  $\gamma$  is sufficiently large the curse of dimensionality is overcome.

In Yang & Barron (1999), it was shown that through the entropy in density spaces defined with distances as in (8), it is possible to determine lower bounds for the learning rate of binary estimators (see Lemma 17). Then, in Petersen & Voigtlaender (2021), a relation between the distances (8) in a specific set of densities and the norms  $L^2(\lambda)$  (being  $\lambda$  the Lebesgue measure on  $[0,1]^d$ ) in  $H_{\mathscr{C}}$  and  $L^1([0,1]^{d-1})$  in  $\mathscr{C}$  was demonstrated, thus adapting the main results of Yang & Barron (1999) to function spaces  $\mathscr{C}$  as in C1). However, the margin condition was not assumed in either Petersen & Voigtlaender (2021) or Yang & Barron (1999).

In this paper, we use an argument to lower bound the minimax expression  $\mathcal{I}_n(\mathscr{C})$  in (4) based on carefully constructing a particular family of probability densities (see Lemma (12)) over a subset of the function class  $\mathscr{C}$ , where the margin condition C3) is satisfied (see Remark 7). Thus, assuming that conditions C1), C2), and C3) hold, our main results provide the following: In Theorem 6, under certain conditions (11), a general lower bound for  $\mathcal{I}_n(\mathscr{C})$  is given when the class  $\mathscr{C}$  has bounded entropies (5) and (6) in the Lebesgue norm  $L^1$ . In Corollary 8, we apply Theorem 6 when the class  $\mathscr{C}$  has entropy bounded in a particular way; then we obtain specific lower bounds for  $\mathcal{I}_n(\mathscr{C})$  that depend on the sample size n, the parameters controlling the entropies, the margin exponent  $\gamma$  in C3) and the  $\alpha$  assumed in the regularity C2) of  $\mathscr{C}$ . Finally, in Corollary 9, we provide explicit lower bounds for  $\mathscr{C}$  being:

• The Barron class,  $\mathscr{C} = \mathcal{B}_C$ , for which we obtain in (13) lower bounds that, for all  $\gamma \geq 3$  (see Remark 10) take the form

$$\mathcal{I}_n(\mathcal{B}_C) \gtrsim \left[ n \left( 1 + 2 \log n \right)^{\frac{(d+1)\gamma}{2(d-1)}} \right]^{-\frac{\gamma}{\gamma + \left( \frac{2(d-1)}{d+1} \right)}},$$

where  $2(d-1)/(d+1) \to 2$  when  $d \to \infty$ . While in García & Petersen (2025) were obtained upper bounds for the learning rate as

$$\mathcal{I}_n(\mathcal{B}_C) \lesssim n^{-\gamma/(2+\gamma)} (1 + \log n), \text{ for all } \gamma > 0.$$

In conclusion, when  $\gamma \geq 3$ , our lower bound matches the upper bound in García & Petersen (2025) up to logarithmic factors. For large dimension d, the exponent of n in both bounds is  $-\gamma/(2+\gamma)$ , showing that the learning rate in García & Petersen (2025) is asymptotically optimal.

• The  $\alpha$ -Hölder continuous class,  $\mathscr{C} = \mathcal{H}_{\alpha}$ , where in (14) we found for all  $\gamma \geq 3\alpha$  (see Remark 10), that

$$\mathcal{I}_n(\mathcal{H}_\alpha) \gtrsim n^{-\frac{1}{1+(d-1)/\gamma}}, \quad \text{for all} \quad \alpha \in (0,1].$$

In contrast, under the noiseless assumption (see Remark 1), an upper bound was obtained in (Kim et al., 2021, Theorem 3.4), as follows

$$\left(\frac{\log^3 n}{n}\right)^{\frac{\alpha}{\alpha+(d-1)/\gamma}}, \quad \text{for all} \quad \gamma \geq 1 \quad \text{and all} \quad \alpha \in (0,1].$$

For values of  $\alpha$  close to 1, the exponents of n in the lower and upper bounds become similar, indicating that the corresponding learning rates are nearly the same, up to logarithmic factors. In particular, for Lipschitz functions ( $\alpha = 1$ ), the learning rate is asymptotically optimal. However, for small values of  $\alpha$ , the lower and upper bounds themselves differ significantly.

• The Convex class,  $\mathscr{C} = \mathcal{C}_B$ , for which in (15), we obtain for all  $\gamma \geq 3$  (see Remark 10),

$$\mathcal{I}_n(\mathcal{C}_B) \gtrsim n^{-\frac{2\gamma}{2\gamma+(d-1)}}$$
.

When  $\gamma$  is large, the lower bound yields learning rates similar to those obtained for general Lipschitz functions in the previous item (case  $\alpha = 1$ ).

#### 2 Main results

In this section, we give the main results of our work. To begin with, we present the following theorem on general bounds for the minimax expression (4).

**Theorem 6.** Let conditions C1), C2) and C3) hold. Let  $V, M: (0, \infty) \to [0, \infty)$  be continuous from the right and non-increasing. If for all  $\varepsilon \in (0,1)$ ,  $\widetilde{\gamma} = \max\{2, \gamma/\alpha - 1\}$  and a constant c > 1, it is satisfied that

$$V_{\mathscr{C},L^{1}}\left(\varepsilon^{\frac{2}{\widetilde{\gamma}+1}}/\left(4\widetilde{\gamma}^{2}(2c)^{\frac{2}{\widetilde{\gamma}+1}}\right)\right)\leq V(\varepsilon)\quad and\quad M_{\mathscr{C},L^{1}}\left((4\widetilde{\gamma})\varepsilon^{\frac{2}{\widetilde{\gamma}+1}}\right)\geq M(\varepsilon). \tag{10}$$

Then, for every pair of sequences  $\{\varepsilon_n\}_{n\in\mathbb{N}}, \{\widetilde{\varepsilon}_n\}_{n\in\mathbb{N}}\subset (0,\infty)$  satisfying

$$n\varepsilon_n^2 = V(\varepsilon_n) = (M(\widetilde{\varepsilon}_n) - 2\log 2)/4,$$
 (11)

we obtain  $\widetilde{\varepsilon}_n^2 \leq (4\varepsilon_n)^2$  and

$$\mathcal{I}_n(\mathscr{C}) \ge \left[ \frac{1}{8(4\widetilde{\gamma})^3} \right] \widetilde{\varepsilon}_n^2, \tag{12}$$

with  $\mathcal{I}_n(\mathscr{C})$  as in (4).

**Remark 7.** For the proof of this theorem, we begin by defining a particular family of densities (18) that satisfy margin condition C3) over the subset  $\mathscr{C}_{\delta} \subseteq \mathscr{C}$  defined in Lemma 12. Then, given that  $\{x_i\}_{i=1}^n$  is distributed according to the margin condition C3) with respect to a measure  $\mu_h$ , in Lemma 13 we define a new class of densities  $\mathcal{P}_{\mathscr{C}_{\delta}}$  over  $\mu$  for which the sample  $S_h = ((x_i, h(x_i)))_{i=1}^n$  is distributed. With this new set of densities  $\mathcal{P}_{\mathscr{C}_{\delta}}$ , we apply Lemma 17 ((Petersen & Voigtlaender, 2021, Theorem 3.11), which is a simplified version of (Yang & Barron, 1999, Theorems 1 and 2)), with a similar argument to the one used in Petersen & Voigtlaender (2021), to obtain our lower bound (12). However, to apply Lemma 17 to the density class  $\mathcal{P}_{\mathscr{C}_{\delta}}$ , it was necessary to prove the results of Lemma 14 and Lemma 15.

As a consequence of the previous theorem, when the entropy of the space  $\mathscr C$  with respect to the Lebesgue  $L^1$  norm is bounded in a particular way, we obtain the following corollary.

**Corollary 8.** Let conditions C1), C2) and C3) hold. Let  $\tilde{\gamma} = \max\{2, \gamma/\alpha - 1\}$ ;  $a \ge 1/2$ ; and  $\mathcal{I}_n(\mathscr{C})$  as in (4). Then, by cases:

I) If 
$$\varepsilon^{-a} \lesssim M_{\mathscr{C},L^1}(\varepsilon) \lesssim \varepsilon^{-a} (1 + \log(1/\varepsilon))$$
, for all  $\varepsilon \in (0,1)$ . Then

$$\mathcal{I}_n(\mathscr{C}) \gtrsim \left[ n \left( 1 + 2 \log n \right)^{\frac{\tilde{\gamma}+1}{a}} \right]^{-\frac{1}{1+\left(\frac{a}{\tilde{\gamma}+1}\right)}}.$$

II) If 
$$V_{\mathscr{C},L^1}(\varepsilon) \approx \varepsilon^{-a}$$
, for all  $\varepsilon \in (0,1)$ . Then

$$\mathcal{I}_n(\mathscr{C}) \gtrsim n^{-\frac{1}{1+\left(\frac{a}{\widetilde{\gamma}+1}\right)}}.$$

Finally, we apply the previous corollary to the spaces defined in Section 1.3 to obtain the next result.

**Corollary 9.** Let  $\gamma > 0$  be the margin exponent in C3). Then, with the notation (4) and the notation in Section 1.3, we obtain:

I) For Barron regular functions,

$$\mathcal{I}_n(\mathcal{B}_C) \gtrsim \left[ n \left( 1 + 2 \log n \right)^{\frac{(d+1) \max\{3,\gamma\}}{2(d-1)}} \right]^{-\frac{\max\{3,\gamma\}}{\max\{3,\gamma\} + \left(\frac{2(d-1)}{d+1}\right)}}.$$

II) For  $\alpha$ -Hölder continuous functions,

for all 
$$\alpha \in (0,1]$$
.

III) For convex functions,

$$\mathcal{I}_n(\mathcal{C}_B) \gtrsim n^{-\frac{2\max\{3,\gamma\}}{2\max\{3,\gamma\}+(d-1)}}.$$
 (15)

(13)

(14)

**Remark 10.** In our results, the term  $\max \{3\alpha, \gamma\}$  appears, which we had to assume since the proof of our main theorem is based on defining the particular family of densities given in (18), for which one of the most important properties is that their square root is Lipschitz in b with respect to b on the  $L^2(\lambda)$  (see (20)).

 $\mathcal{I}_n(\mathcal{H}_\alpha) \gtrsim n^{-\frac{\max\{3\alpha,\gamma\}}{\max\{3\alpha,\gamma\}+(d-1)}},$ 

In Kim et al. (2021) and García & Petersen (2025), upper bounds on the learning rate of approximating, by deep neural networks, binary classifiers with decision boundary satisfying the regularity condition C1) in function classes  $\mathcal{H}_{\alpha}$  and  $\mathcal{B}_{\mathcal{C}}$ , respectively, were shown. See discussion of our results in comparison with those of Kim et al. (2021) and García & Petersen (2025) in Section 1.4.

**Remark 11.** The main results Theorem 6, Corollary 8 and Corollary 9 identify precise optimal values for the learning rates of a binary classifier from noiseless samples, if a margin condition is satisfied and the decision boundary stems from a general class of functions. There are considerable limitations of these results that we would like to stress:

1. The established rates are understood in the inf-sup sense where all distributions that satisfy the margin condition are considered. In practice, many more favorable conditions on the data distribution may hold, which shows that the observed learning rates could be considerably faster.

2. As mentioned in Remark 1 we consider here noiseless data. This assumption is crucial to identify meaningful lower bounds, since it is otherwise unclear if the lower bounds stem from the noise or from the properties of the decision boundary (see (Petersen & Voigtlaender, 2021, Section 1.1, Point 1)). Still, noisy data appears extremely often in applications, and that situation is not covered by our main results.

3. Our approach requires precise bounds on the covering numbers of the spaces that model the decision boundary. We have described three such cases in this manuscript. However, for many, more exotic spaces, such estimates may not exist.

# REFERENCES

A. R. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information Theory*, 39(3):930–945, 1993. doi: 10.1109/18.256500.

E. M. Bronshtein.  $\varepsilon$ -entropy of convex sets and functions. *Siberian Mathematical Journal*, 17: 393–398, 1976. doi: 10.1007/BF00967858.

Andrei Caragea, Philipp Petersen, and Felix Voigtlaender. Neural network approximation and estimation of classifiers with classification boundary in a Barron class. *The Annals of Applied Probability*, 33(4):3039 – 3079, 2023. doi: 10.1214/22-AAP1884.

Andreas Christmann and Ingo Steinwart. *Support Vector Machines*. Springer-Verlag New York, 2008. doi: 10.1007/978-0-387-77242-4.

G. F. Clements. Entropies of several sets of real valued functions. *Pacific Journal of Mathematics*, 13, 1963. doi: 10.2140/pjm.1963.13.1085.

- D. Dryanov. Kolmogorov entropy for classes of convex functions. *Constructive Approximation*, 30: 137–153, 2009. doi: 10.1007/s00365-009-9053-3.
  - Jonathan García and Philipp Petersen. High-dimensional classification problems with barron regular boundaries under margin conditions. *arXiv:2412.07312*, 2025. doi: 10.48550/arXiv.2412.07312.
  - Adityanand Guntuboyina and Bodhisattva Sen. Covering numbers for convex functions. *IEEE Transactions on Information Theory*, 59(4):1957–1965, 2013. doi: 10.1109/TIT.2012.2235172.
  - Witold Hurewicz and Henry Wallman. *Dimension theory*. Princeton University Press, 1948. URL https://mathscinet.ams.org/mathscinet-getitem?mr=6493.
  - Yongdai Kim, Ilsang Ohn, and Dongha Kim. Fast convergence rates of deep neural networks for classification. *Neural Networks*, 138:179–197, 2021. doi: 10.1016/j.neunet.2021.02.012.
  - A. N. Kolmogorov and V. M. Tihomirov.  $\varepsilon$ -entropy and  $\varepsilon$ -capacity of sets in functional spaces. *Springer Nature*, pp. 86–170, 1993. doi: 10.1007/978-94-017-2973-4\_7.
  - David Krieg and Mathias Sonnleitner. Random points are optimal for the approximation of sobolev functions. *IMA Journal of Numerical Analysis*, 44(3):1346–1371, 2023. doi: 10.1093/imanum/drad014.
  - Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 2009. URL https://www.cs.toronto.edu/~kriz/cifar.html.
  - Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. doi: 10.1109/5.726791.
  - Philipp Petersen and Felix Voigtlaender. Optimal learning of high-dimensional classification problems using deep neural networks. *arXiv:2112.12555*, 2021. doi: 10.48550/arXiv.2112.12555.
  - Charles J. Stone. Optimal Global Rates of Convergence for Nonparametric Regression. *The Annals of Statistics*, 10(4):1040 1053, 1982. doi: 10.1214/aos/1176345969.
  - Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv:1708.07747*, 2017. doi: 10.48550/arXiv.1708.07747.
  - Yuhong Yang and Andrew Barron. Information-theoretic determination of minimax rates of convergence. *The Annals of Statistics*, 27(5):1564 1599, 1999. doi: 10.1214/aos/1017939142.

## A AUXILIARY RESULTS

In this appendix, we present our auxiliary results used to prove Theorem 6.

**Lemma 12.** Let  $\mathscr{C}$  satisfy C2) with  $\alpha \in (0,1]$  and K > 0. Then, the following statements hold.

1) For all  $x \in [0,1]^d$ ,

$$\left(\frac{1}{\widetilde{K}}|b(\boldsymbol{x}^{(d)}) - x_d|\right)^{\frac{1}{\alpha}} \le \operatorname{dist}(\boldsymbol{x}, \partial\Omega_{h_b}) \le |b(\boldsymbol{x}^{(d)}) - x_d|, \tag{16}$$

where

$$1 < \widetilde{K} := \widetilde{K}(\alpha) = \begin{cases} \sqrt{1 + K^2} & \text{if } \alpha = 1\\ \max\{2^{\alpha}, 2K\} & \text{if } \alpha < 1. \end{cases}$$
 (17)

2) Let  $\delta \in (0, (4\widetilde{K})^{-\frac{1}{\alpha}}]$  and  $\mathscr{C}_{\delta} := \{\delta b : b \in \mathscr{C}\}$ . Let  $\widetilde{f}_{h_b} : [0, 1]^d \to [0, \infty)$  be defined by

$$\widetilde{f}_{h_b}(\boldsymbol{x}) := \begin{cases} |b(\boldsymbol{x}^{(d)}) - x_d|^{\widetilde{\gamma}} & \text{if } x_d \le 1/2\\ \widetilde{C}_b & \text{if } x_d > 1/2, \end{cases}$$
(18)

with  $\widetilde{\gamma} := \max\left\{2, \frac{\gamma}{\alpha} - 1\right\}, \gamma > 0$  and

434
435
436
$$\widetilde{C}_b := 2 \left( 1 - \int_{[0,1]^{d-1} \times [0,1/2]} |b(\boldsymbol{x}^{(d)}) - x_d|^{\widetilde{\gamma}} d\boldsymbol{x} \right). \tag{19}$$

Then,  $\mathcal{C}_{\delta} \subseteq \mathcal{C}$  and for all  $b \in \mathcal{C}_{\delta}$ , we obtain that  $\widetilde{f}_{h_b}$  is a Lebesgue density inducing a measure  $\widetilde{\mu}_{h_b}$  that satisfies the margin condition C3), with margin exponent  $\gamma$  and for all  $0 < \epsilon \leq \delta$ .

Moreover,  $\sqrt{\tilde{f}_{h_b}}$  is Lipschitz in b with respect to b on the  $L^2(\lambda)$  norm, namely

$$\left\| \sqrt{\widetilde{f}_{h_{b_1}}} - \sqrt{\widetilde{f}_{h_{b_2}}} \right\|_{L^2(\lambda)} \le (3/8)\widetilde{\gamma} \|b_1 - b_2\|_{L^2([0,1]^{d-1})}, \quad \text{for all} \quad b_1, b_2 \in \mathscr{C}_{\delta}. \tag{20}$$

Proof. By cases:

 1) For all  $x' \in \partial \Omega_{h_h} = \{ x \in [0,1]^d : b(x^{(d)}) = x_d \}$ , and any  $x \in [0,1]^d$ ,

$$\|\boldsymbol{x} - \boldsymbol{x'}\|^2 = \|\boldsymbol{x}^{(d)} - \boldsymbol{x'}^{(d)}\|^2 + |b(\boldsymbol{x'}^{(d)}) - x_d|^2.$$
 (21)

Let  $t := \| \boldsymbol{x}^{(d)} - \boldsymbol{x'}^{(d)} \|$  and  $w := |b(\boldsymbol{x}^{(d)}) - x_d|$ . If  $w - Kt^{\alpha} \le 0$ , we get  $t^2 \ge (w/K)^{2/\alpha}$ , and

$$\|\boldsymbol{x} - \boldsymbol{x'}\|^2 = t^2 + |b(\boldsymbol{x'}^{(d)}) - x_d|^2 \ge (w/K)^{2/\alpha}.$$
 (22)

If  $w - Kt^{\alpha} \ge 0$ , condition C2) implies

$$w = |b(\boldsymbol{x}^{(d)}) - x_d| \le |b(\boldsymbol{x'}^{(d)}) - x_d| + |b(\boldsymbol{x}^{(d)}) - b(\boldsymbol{x'}^{(d)})$$

$$\le |b(\boldsymbol{x'}^{(d)}) - x_d| + K \|\boldsymbol{x}^{(d)} - \boldsymbol{x'}^{(d)}\|^{\alpha}$$

$$= |b(\boldsymbol{x'}^{(d)}) - x_d| + Kt^{\alpha},$$

therefore

$$\|\boldsymbol{x} - \boldsymbol{x'}\|^2 = t^2 + |b(\boldsymbol{x'}^{(d)}) - x_d|^2 \ge t^2 + (w - Kt^{\alpha})^2.$$
 (23)

So, when  $\alpha = 1$ , it follows that

$$\|\boldsymbol{x} - \boldsymbol{x'}\|^{2} \ge (1 + K^{2})t^{2} - 2wKt + w^{2}$$

$$\ge \frac{4w^{2}(1 + K^{2}) - (2wK)^{2}}{4(1 + K^{2})}$$

$$= \frac{w^{2}}{1 + K^{2}}.$$
(24)

Otherwise,  $\alpha < 1$  and we set a threshold  $t_0 := (w/(2K))^{1/\alpha}$  to minimize the right hand side of (23) on t. By cases:

• If  $t \le t_0$ , we know that  $w - Kt^{\alpha} \ge w - Kt^{\alpha}_0 = w/2 > 0$ , therefore

$$\|\boldsymbol{x} - \boldsymbol{x'}\|^2 \ge t^2 + (w - Kt_0^{\alpha})^2 = t^2 + (w/2)^2 \ge (w/2)^2.$$
 (25)

• If  $t > t_0$ , we obtain

$$\|\boldsymbol{x} - \boldsymbol{x'}\|^2 \ge t^2 + (w - Kt^{\alpha})^2 \ge t^2 > (w/(2K))^{2/\alpha}.$$
 (26)

Thus, (22), (24), (25) and (26), imply that

$$\|oldsymbol{x} - oldsymbol{x'}\| \ge \left(rac{1}{\widetilde{K}}|b(oldsymbol{x}^{(d)}) - x_d|
ight)^{rac{1}{lpha}},$$

with  $\widetilde{K}$  as in (17). Then, by the above inequality and identity (21), we arrive at

$$\left(\frac{1}{\widetilde{K}}|b(\boldsymbol{x}^{(d)}) - x_d|\right)^{\frac{1}{\alpha}} \leq \inf_{\boldsymbol{x'} \in \partial \Omega_{h_h}} \|\boldsymbol{x} - \boldsymbol{x'}\| \leq |b(\boldsymbol{x}^{(d)}) - x_d|,$$

where the upper bound was obtained given that  $(x_1, \ldots, x_{d-1}, b(\mathbf{x}^{(d)})) \in \partial \Omega_{h_b}$ .

2) It follows from C2) that  $\delta b \in \mathscr{C}$  for all  $b \in \mathscr{C}$ , since  $\delta < 1$ . So,  $\mathscr{C}_{\delta} \subseteq \mathscr{C}$ .

For any  $b\in\mathscr{C}_{\delta}$ , to show that  $\widetilde{f}_{h_b}$  is a density function with respect to  $\lambda$ , we get

$$\begin{split} \int_{[0,1]^d} \widetilde{f}_{h_b} d\lambda &= \int_{[0,1]^{d-1} \times [0,1/2]} \widetilde{f}_{h_b} d\lambda + \int_{[0,1]^{d-1} \times [1/2,1]} \widetilde{f}_{h_b} d\lambda \\ &= \int_{[0,1]^{d-1} \times [0,1/2]} |b(\boldsymbol{x}^{(d)}) - x_d|^{\widetilde{\gamma}} d\boldsymbol{x} + \int_{[0,1]^{d-1} \times [1/2,1]} \widetilde{C}_b d\boldsymbol{x} \\ &= \int_{[0,1]^{d-1} \times [0,1/2]} |b(\boldsymbol{x}^{(d)}) - x_d|^{\widetilde{\gamma}} d\boldsymbol{x} + (1/2) \widetilde{C}_b \\ &= 1. \end{split}$$

Now, for all  $b \in \mathscr{C}_{\delta}$ , we see that  $\widetilde{f}_{h_b}$  satisfies the margin condition C3) for all  $0 < \epsilon \le \delta$ , as follows. Using (16), we obtain that

$$|x_d - b(x^{(d)})| \le |b(x^{(d)}) - x_d| \le \widetilde{K} \left( \operatorname{dist}(x, \partial \Omega_{h_b}) \right)^{\alpha} \le \widetilde{K} \epsilon^{\alpha} \le \widetilde{K} \delta^{\alpha}, \quad \text{for all} \quad x \in B_{\epsilon}^{h_b}.$$

Therefore.

$$x_d \leq \widetilde{K}\delta^{\alpha} + b(\boldsymbol{x}^{(d)}) \leq \widetilde{K}\delta^{\alpha} + \delta \leq 2\widetilde{K}\delta^{\alpha} \leq 2\widetilde{K}/(4\widetilde{K}) = 1/2$$
 for all  $\boldsymbol{x} \in B_{\epsilon}^{h_b}$ .

Then,  $B_{\epsilon}^{h_b} \subseteq [0,1]^{d-1} \times [0,1/2]$ , and

$$\begin{split} \widetilde{\mu}_{h_b}(B^{h_b}_{\epsilon}) &= \widetilde{\mu}_{h_b}(B^{h_b}_{\epsilon} \cap ([0,1]^{d-1} \times [0,1/2])) \\ &= \int_{B^{h_b}_{\epsilon}} |b(\boldsymbol{x}^{(d)}) - x_d|^{\widetilde{\gamma}} d\lambda \\ &\leq \int_{\left\{\boldsymbol{x} \in [0,1]^d: |b(\boldsymbol{x}^{(d)}) - x_d| \leq \widetilde{K} \epsilon^{\alpha}\right\}} |b(\boldsymbol{x}^{(d)}) - x_d|^{\widetilde{\gamma}} d\lambda \\ &\leq (\widetilde{K} \epsilon^{\alpha})^{\widetilde{\gamma}} \int_{\left\{\boldsymbol{x} \in [0,1]^d: |b(\boldsymbol{x}^{(d)}) - x_d| \leq \widetilde{K} \epsilon^{\alpha}\right\}} 1 d\lambda \\ &\leq \widetilde{K}^{\widetilde{\gamma}} \epsilon^{\max\{2\alpha, \gamma - \alpha\}} (2\widetilde{K} \epsilon^{\alpha}) \\ &\leq 2\widetilde{K}^{\widetilde{\gamma} + 1} \epsilon^{\max\{3\alpha, \gamma\}} \\ &\leq \left(2\widetilde{K}^{\widetilde{\gamma} + 1}\right) \epsilon^{\gamma}. \end{split}$$

This proves that condition C3) is satisfied.

Moreover, as a consequence of the mean value theorem, it is well known that

$$||u|^{\beta} - |v|^{\beta}| \le \beta \max\{|u|, |v|\}^{\beta - 1} ||u| - |v||$$

$$\le \beta |u - v| \quad \text{for all} \quad u, v \in [-1, 1] \quad \text{and} \quad \beta \ge 1.$$
(27)

So, we apply (27) by cases:

• Let  $g_b(\boldsymbol{x}) := |b(\boldsymbol{x}^{(d)}) - x_d|^{\widetilde{\gamma}/2}$ . Then, for all  $b_1, b_2 \in \mathscr{C}_{\delta}$  and all  $\boldsymbol{x} \in [0, 1]^d$ ,

$$|g_{b_1}(\boldsymbol{x}) - g_{b_2}(\boldsymbol{x})| = \left| |b_1(\boldsymbol{x}^{(d)}) - x_d|^{\widetilde{\gamma}/2} - |b_2(\boldsymbol{x}^{(d)}) - x_d|^{\widetilde{\gamma}/2} \right|$$
  
 
$$\leq (\widetilde{\gamma}/2)|b_1(\boldsymbol{x}^{(d)}) - b_2(\boldsymbol{x}^{(d)})|.$$

Therefore,

$$\int_{[0,1]^{d-1}\times[0,1/2]} |g_{b_1}(\boldsymbol{x}) - g_{b_2}(\boldsymbol{x})|^2 d\boldsymbol{x} \le (\widetilde{\gamma}/2)^2 \int_{[0,1]^{d-1}\times[0,1/2]} |b_1(\boldsymbol{x}^{(d)}) - b_2(\boldsymbol{x}^{(d)})|^2 d\boldsymbol{x} 
= (\widetilde{\gamma}^2/8) \|b_1 - b_2\|_{L^2([0,1]^{d-1})}^2.$$

Thus,  $g_b$  is Lipschitz with respect to b on their respective  $L^2$  norms.

• For all  $b_1, b_2 \in \mathscr{C}_{\delta}$ ,

$$\frac{1}{4} \left| \widetilde{C}_{b_{1}} - \widetilde{C}_{b_{2}} \right|^{2} = \left( \int_{[0,1]^{d-1} \times [0,1/2]} \left( g_{b_{2}}^{2}(\boldsymbol{x}) - g_{b_{1}}^{2}(\boldsymbol{x}) \right) d\boldsymbol{x} \right)^{2} \\
\leq \int_{[0,1]^{d-1} \times [0,1/2]} \left( g_{b_{2}}(\boldsymbol{x}) - g_{b_{1}}(\boldsymbol{x}) \right)^{2} \left( g_{b_{2}}(\boldsymbol{x}) + g_{b_{1}}(\boldsymbol{x}) \right)^{2} d\boldsymbol{x} \\
\leq 4 \int_{[0,1]^{d-1} \times [0,1/2]} \left| g_{b_{2}}(\boldsymbol{x}) - g_{b_{1}}(\boldsymbol{x}) \right|^{2} d\boldsymbol{x} \\
\leq (\widetilde{\gamma}^{2}/2) \left\| b_{1} - b_{2} \right\|_{L^{2}([0,1]^{d-1})}^{2}.$$

Also,

$$\left(\sqrt{\widetilde{C}_{b_1}} - \sqrt{\widetilde{C}_{b_2}}\right)^2 \le \frac{1}{4} \left(\sqrt{\widetilde{C}_{b_1}} + \sqrt{\widetilde{C}_{b_2}}\right)^2 \left(\sqrt{\widetilde{C}_{b_1}} - \sqrt{\widetilde{C}_{b_2}}\right)^2$$
$$= \frac{1}{4} \left|\widetilde{C}_{b_1} - \widetilde{C}_{b_2}\right|^2,$$

since for all  $i \in \{1, 2\}$ ,

$$\widetilde{C}_{b_i} \ge 2\left(1 - \int_{[0,1]^{d-1} \times [0,1/2]} 1d\mathbf{x}\right) = 1.$$
 (28)

Then,

$$\left(\sqrt{\widetilde{C}_{b_1}} - \sqrt{\widetilde{C}_{b_2}}\right)^2 \le (\widetilde{\gamma}^2/2) \|b_1 - b_2\|_{L^2([0,1]^{d-1})}^2.$$

Using the items above, we conclude that

$$\begin{split} & \left\| \sqrt{\widetilde{f}_{h_{b_{1}}}} - \sqrt{\widetilde{f}_{h_{b_{2}}}} \right\|_{L^{2}(\lambda)}^{2} \\ &= \int_{[0,1]^{d-1} \times [0,1/2]} \left| g_{b_{1}}(\boldsymbol{x}) - g_{b_{2}}(\boldsymbol{x}) \right|^{2} d\boldsymbol{x} + \int_{[0,1]^{d-1} \times [1/2,1]} \left( \sqrt{\widetilde{C}_{b_{1}}} - \sqrt{\widetilde{C}_{b_{2}}} \right)^{2} d\boldsymbol{x} \\ &\leq (\widetilde{\gamma}^{2}/8) \left\| b_{1} - b_{2} \right\|_{L^{2}([0,1]^{d-1})}^{2} + (\widetilde{\gamma}^{2}/4) \left\| b_{1} - b_{2} \right\|_{L^{2}([0,1]^{d-1})}^{2} \\ &\leq (3/8)\widetilde{\gamma}^{2} \left\| b_{1} - b_{2} \right\|_{L^{2}([0,1]^{d-1})}^{2}, \quad \text{for all} \quad b_{1}, b_{2} \in \mathscr{C}_{\delta}. \end{split}$$

Therefore (20) is fulfilled.

**Lemma 13.** Let  $\mathcal{P}_{\mathscr{C}_{\delta}} := \{p_h : h \in H_{\mathscr{C}_{\delta}}\}$ , where

$$p_h(\boldsymbol{x}, y) := \begin{cases} h(\boldsymbol{x}) \widetilde{f}_h(\boldsymbol{x}) / \nu_1 & \text{if } y = 1\\ (1 - h(\boldsymbol{x})) \widetilde{f}_h(\boldsymbol{x}) / \nu_0 & \text{if } y = 0, \end{cases}$$
(29)

and  $\tilde{f}_h$  is a Lebesgue density that defines a measure  $\tilde{\mu}_h$  as in (18). Then,  $p_h$  is a density function with respect to the measure  $\mu$  as in (3). Furthermore,

$$\mathbf{x} \sim \widetilde{\mu}_h \text{ if only if } (\mathbf{x}, h(\mathbf{x})) \sim p_h \text{ (with respect to } \boldsymbol{\mu}).$$
 (30)

*Proof.* For all  $h \in H_{\mathscr{C}_{\delta}}$ , we know that  $p_h \in D_{\mu}$ , since

$$\int_{\Lambda} p_h(\boldsymbol{x}, y) d\boldsymbol{\mu} = \nu_0 \int_{[0,1]^d} p_h(\boldsymbol{x}, 0) d\lambda + \nu_1 \int_{[0,1]^d} p_h(\boldsymbol{x}, 1) d\lambda$$

$$= \int_{[0,1]^d} (1 - h(\boldsymbol{x})) d\widetilde{\mu}_h + \int_{[0,1]^d} h(\boldsymbol{x}) d\widetilde{\mu}_h$$

$$= \int_{[0,1]^d} 1 d\widetilde{\mu}_h = 1.$$

In addition, for  $(x, h(x)) \in M \subset \Lambda$  and  $x \sim \widetilde{\mu}_h$ , we have that

$$\begin{split} \int_{[0,1]^d} \mathbb{1}_M(\boldsymbol{x}, h(\boldsymbol{x})) d\widetilde{\mu}_h &= \int_{[0,1]^d} \mathbb{1}_M(\boldsymbol{x}, \mathbb{1}_{\Omega_h}(\boldsymbol{x})) d\widetilde{\mu}_h \\ &= \int_{[0,1]^d \setminus \Omega_h} \mathbb{1}_M(\boldsymbol{x}, 0) d\widetilde{\mu}_h + \int_{\Omega_h} \mathbb{1}_M(\boldsymbol{x}, 1) d\widetilde{\mu}_h \\ &= \int_{[0,1]^d} \mathbb{1}_M(\boldsymbol{x}, 0) \left(1 - \mathbb{1}_{\Omega_h}(\boldsymbol{x})\right) d\widetilde{\mu}_h + \int_{[0,1]^d} \mathbb{1}_M(\boldsymbol{x}, 1) \mathbb{1}_{\Omega_h}(\boldsymbol{x}) d\widetilde{\mu}_h \\ &= \int_{[0,1]^d} \mathbb{1}_M(\boldsymbol{x}, 0) \left(1 - h(\boldsymbol{x})\right) d\widetilde{\mu}_h + \int_{[0,1]^d} \mathbb{1}_M(\boldsymbol{x}, 1) h(\boldsymbol{x}) d\widetilde{\mu}_h \\ &= \nu_0 \int_{[0,1]^d} \mathbb{1}_M(\boldsymbol{x}, 0) p_h(\boldsymbol{x}, 0) d\lambda + \nu_1 \int_{[0,1]^d} \mathbb{1}_M(\boldsymbol{x}, 1) p_h(\boldsymbol{x}, 1) d\lambda \\ &= \int_{\Lambda} \mathbb{1}_M(\boldsymbol{x}, y) p_h(\boldsymbol{x}, y) d\boldsymbol{\mu}, \end{split}$$

equivalent to  $(x, h(x)) \sim p_h$  with respect to  $\mu$  and (30) is satisfied.

**Lemma 14.** Let  $h_1, h_2 \in H_{\mathscr{C}_{\delta}}$  be defined by  $b_1, b_2 \in \mathscr{C}_{\delta}$ , respectively. Then,

$$\|h_1 - h_2\|_{L^2(\widetilde{\mu}_{h_i})}^2 = \frac{1}{\widetilde{\gamma} + 1} \|b_1 - b_2\|_{L^{\widetilde{\gamma} + 1}([0, 1]^{d - 1})}^{\widetilde{\gamma} + 1} \quad \text{for all} \quad i \in \{1, 2\},$$
 (31)

with  $\widetilde{\mu}_{h_i}$  and  $\widetilde{\gamma}$  as in Lemma 12. Moreover, for all  $p_{h_1}, p_{h_2} \in \mathcal{P}_{\mathscr{C}_{\delta}}$ , we have that

$$\|h_1 - h_2\|_{L^2(\widetilde{\mu}_{h_i})}^2 \le \rho_H^2(p_{h_1}, p_{h_2}) \le 2\widetilde{\gamma}^2 \|b_1 - b_2\|_{L^1([0,1]^{d-1})}, \tag{32}$$

for all  $i \in \{1, 2\}$ .

*Proof.* For all  $i \in \{1, 2\}$ , we know that  $b_i \in \mathscr{C}_{\delta}$  and therefore  $b_i(\boldsymbol{x}^{(d)}) \leq \delta$  for all  $\boldsymbol{x}^{(d)} \in [0, 1]^{d-1}$ . Also,  $\delta \in (0, (4\widetilde{K})^{-\frac{1}{\alpha}}]$ , so  $\delta \leq (4\widetilde{K})^{-1} < 1/2$ . Then, by definition (1),

$$G_{1,2} := \left\{ \boldsymbol{x} \in [0,1]^d : h_1(\boldsymbol{x}) \neq h_2(\boldsymbol{x}) \right\}$$

$$= \left\{ \boldsymbol{x} \in [0,1]^d : h_1(\boldsymbol{x}) = 0 \land h_2(\boldsymbol{x}) = 1 \right\} \cup \left\{ \boldsymbol{x} \in [0,1]^d : h_1(\boldsymbol{x}) = 1 \land h_2(\boldsymbol{x}) = 0 \right\}$$

$$= \left\{ \boldsymbol{x} \in [0,1]^d : b_2(\boldsymbol{x}^{(d)}) \leq x_d < b_1(\boldsymbol{x}^{(d)}) \right\} \cup \left\{ \boldsymbol{x} \in [0,1]^d : b_1(\boldsymbol{x}^{(d)}) \leq x_d < b_2(\boldsymbol{x}^{(d)}) \right\}$$

$$= \left\{ \boldsymbol{x} \in [0,1]^d : \min_{i \in \{1,2\}} \left\{ b_i(\boldsymbol{x}^{(d)}) \right\} \leq x_d < \max_{i \in \{1,2\}} \left\{ b_i(\boldsymbol{x}^{(d)}) \right\} < 1/2 \right\}$$

$$\subseteq [0,1]^{d-1} \times [0,1/2].$$
(33)

Furthermore,  $h_1$  and  $h_2$  are indicator functions, thus  $|h_1 - h_2| = \mathbb{1}_{G_{1,2}}$ , and using (33), we get

$$\begin{aligned} \|h_1 - h_2\|_{L^2(\widetilde{\mu}_{h_i})}^2 &= \int_{[0,1]^d} |h_1 - h_2|^2 \widetilde{f}_{h_i} d\lambda \\ &= \int_{G_{1,2}} |b_i(\boldsymbol{x}^{(d)}) - x_d|^{\widetilde{\gamma}} d\boldsymbol{x} \\ &= \frac{1}{\widetilde{\gamma} + 1} \int_{[0,1]^{d-1}} |b_1(\boldsymbol{x}^{(d)}) - b_2(\boldsymbol{x}^{(d)})|^{\widetilde{\gamma} + 1} d\boldsymbol{x}^{(d)} \\ &= \frac{1}{\widetilde{\gamma} + 1} \|b_1 - b_2\|_{L^{\widetilde{\gamma} + 1}([0,1]^{d-1})}^{\widetilde{\gamma} + 1} \quad \text{for all} \quad i \in \{1, 2\}. \end{aligned}$$

From the above identity, it follows (31). Moreover,  $h_i = \sqrt{h_i}$  and  $1 - h_i = \sqrt{1 - h_i}$ , for all  $i \in \{1, 2\}$ . Hence,

$$\rho_H^2(p_{h_1}, p_{h_2}) - 2 = -2 \int_{\Lambda} \sqrt{p_{h_1} p_{h_2}} d\mu$$

$$= -2 \left[ \int_{[0,1]^d \times \{1\}} \sqrt{p_{h_1} p_{h_2}} d\mu + \int_{[0,1]^d \times \{0\}} \sqrt{p_{h_1} p_{h_2}} d\mu \right]$$

$$= -2 \left[ \int_{[0,1]^d} \sqrt{\widetilde{f}_{h_1} \widetilde{f}_{h_2}} \left( \sqrt{h_1 h_2} + \sqrt{(1 - h_1)(1 - h_2)} \right) d\lambda \right]$$

$$= 2 \left[ \int_{[0,1]^d} \sqrt{\widetilde{f}_{h_1} \widetilde{f}_{h_2}} \left( (\sqrt{h_1} - 2h_1 h_2 + h_2 - 1) d\lambda \right) \right]$$

$$= 2 \left[ \int_{[0,1]^d} \sqrt{\widetilde{f}_{h_1} \widetilde{f}_{h_2}} \left( (\sqrt{h_1} - \sqrt{h_2})^2 - 1 \right) d\lambda \right]$$

$$= 2 \left[ \int_{[0,1]^d} \sqrt{\widetilde{f}_{h_1} \widetilde{f}_{h_2}} \left( (h_1 - h_2 - 1) d\lambda \right) \right].$$

We see that

$$\begin{split} 2 - 2 \int_{[0,1]^d} \sqrt{\widetilde{f}_{h_1} \widetilde{f}_{h_2}} \, d\lambda &= \int_{[0,1]^d} \widetilde{f}_{h_1} d\lambda - 2 \int_{[0,1]^d} \sqrt{\widetilde{f}_{h_1} \widetilde{f}_{h_2}} \, d\lambda + \int_{[0,1]^d} \widetilde{f}_{h_2} d\lambda \\ &= \int_{[0,1]^d} \left( \sqrt{\widetilde{f}_{h_1}} - \sqrt{\widetilde{f}_{h_2}} \right)^2 d\lambda. \end{split}$$

Then the above two identities imply

$$\rho_H^2(p_{h_1}, p_{h_2}) = 2 \int_{[0,1]^d} \sqrt{\widetilde{f}_{h_1} \widetilde{f}_{h_2}} |h_1 - h_2| \, d\lambda + \int_{[0,1]^d} \left( \sqrt{\widetilde{f}_{h_1}} - \sqrt{\widetilde{f}_{h_2}} \right)^2 d\lambda. \tag{34}$$

Now we bound show (32) by cases:

• Upper bound: For all  $i \in \{1,2\}$ , we know that  $\widetilde{f}_{h_i} \leq 3$ , since (19) and (28) implies

$$|b_i(\boldsymbol{x}^{(d)}) - x_d|^{\widetilde{\gamma}} \le 1 \le C_b \le 2 \left(1 + \int_{[0,1]^{d-1} \times [0,1/2]} |b_i(\boldsymbol{x}^{(d)}) - x_d|^{\widetilde{\gamma}} d\boldsymbol{x}\right) \le 3.$$

We use (20), (33) and (34), to obtain

$$\rho_H^2(p_{h_1}, p_{h_2}) \le 6 \int_{[0,1]^d} |h_1 - h_2| \, d\lambda + \left\| \sqrt{\widetilde{f}_{h_1}} - \sqrt{\widetilde{f}_{h_2}} \right\|_{L^2(\lambda)}^2 \\
\le 6 \int_{G_{1,2}} 1 \, d\lambda + \widetilde{\gamma}^2 \|b_1 - b_2\|_{L^2([0,1]^{d-1})}^2 \\
= 6 \|b_1 - b_2\|_{L^1([0,1]^{d-1})} + (3/8)\widetilde{\gamma}^2 \|b_1 - b_2\|_{L^2([0,1]^{d-1})}^2,$$

and we know that  $|b_1(z) - b_2(z)| \le 1$  for all  $z \in [0, 1]^{d-1}$ , then

$$\rho_H^2(p_{h_1}, p_{h_2}) \le (6 + (3/8)\widetilde{\gamma}^2) \|b_1 - b_2\|_{L^1([0,1]^{d-1})} \le 2\widetilde{\gamma}^2 \|b_1 - b_2\|_{L^1([0,1]^{d-1})}.$$

• Lower bound: By (33), (34) and as  $|h_1 - h_2| = \mathbb{1}_{G_{1,2}}$ , we get

$$\begin{split} \rho_H^2(p_{h_1},p_{h_2}) &= 2 \int_{G_{1,2}} \sqrt{\widetilde{f}_{h_1}\widetilde{f}_{h_2}} \, d\lambda + \int_{([0,1]^d \backslash G_{1,2}) \cup G_{1,2}} \left( \sqrt{\widetilde{f}_{h_1}} - \sqrt{\widetilde{f}_{h_2}} \right)^2 d\lambda \\ &= \int_{G_{1,2}} (\widetilde{f}_{h_1} + \widetilde{f}_{h_2}) d\lambda + \int_{([0,1]^d \backslash G_{1,2})} \left( \sqrt{\widetilde{f}_{h_1}} - \sqrt{\widetilde{f}_{h_2}} \right)^2 d\lambda \\ &\geq \int_{[0,1]^d} |h_1 - h_2|^2 \widetilde{f}_{h_i} d\lambda \\ &= \|h_1 - h_2\|_{L^2(\widetilde{\mu}_{h_i})}^2 \,, \quad \text{for all} \quad i \in \{1,2\}. \end{split}$$

**Lemma 15.** Under the conditions C1), C2) and C3), let  $\mathscr{C}(\xi) \subseteq \mathscr{C}_{\delta} \subseteq \mathscr{C}$  be a  $\xi$ -packing set as in (P), with respect to the  $L^1$  norm, and  $\xi \in (0,1)$ . Let

$$\mathscr{C}^{\xi} := \left\{ \xi^{\widetilde{\gamma}} b : b \in \mathscr{C}(\xi) \right\} \subseteq \mathscr{C}(\xi) \quad \text{and} \quad \mathcal{P}_{\mathscr{C}^{\xi}} := \left\{ p_{h_b} \in \mathcal{P}_{\mathscr{C}_{\delta}} : b \in \mathscr{C}^{\xi} \right\},$$

with  $\widetilde{\gamma} = \max\{2, \gamma/\alpha - 1\}$  as in Lemma 12. Then,

$$\mathcal{I}_{n}(\mathscr{C}) \geq \frac{1}{(4\widetilde{\gamma})^{3}} \left[ \inf_{E \in \mathcal{A}_{n}(\mathcal{D}_{\mu})} \sup_{p' \in \mathcal{P}_{\mathscr{C}} \xi} \mathbb{E}_{\mathbf{S}} \stackrel{iid}{\sim} p' \rho_{H}^{2}(E(\mathbf{S}), p') \right], \tag{35}$$

where  $\mathcal{I}_n(\mathscr{C})$  is defined in (4). Moreover,

$$V_{\mathcal{P}_{\mathscr{C}^{\xi}},\rho_{H}}(\varepsilon) \leq V_{\mathscr{C},L^{1}}\left(\varepsilon^{\frac{2}{\gamma+1}}/(4\widetilde{\gamma}^{2})\right) \quad and \quad M_{\mathcal{P}_{\mathscr{C}^{\xi}},\rho_{H}}(\varepsilon) \geq M_{\mathscr{C},L^{1}}\left((4\widetilde{\gamma})\varepsilon^{\frac{2}{\gamma+1}}\right), \tag{36}$$

for all  $\varepsilon > 0$ .

*Proof.* We begin by bounding  $\mathcal{I}_n(\mathscr{C})$  from below. In order to show that

$$4\mathcal{I}_{n}(\mathscr{C}) \geq \widetilde{\mathcal{I}_{n}}(\mathscr{C}_{\delta}) := \inf_{A \in \mathcal{A}_{n}(H_{\mathscr{C}_{\delta}})} \sup_{\substack{h \in H_{\mathscr{C}_{\delta}} \\ \mu_{h} \text{ has density as in (18)}}} \mathbb{E}_{\{\boldsymbol{x}_{i}\}_{i=1}^{n} \stackrel{iid}{\sim} \mu_{h}} \|A(\boldsymbol{S}_{h}) - h\|_{L^{2}(\mu_{h})}^{2}$$
(37)

is satisfied, we introduce (Petersen & Voigtlaender, 2021, Lemma 3.2).

**Lemma 16.** Let  $(X, \rho)$  be a metric space, with distance function  $\rho$ , and let  $\emptyset \neq M \subset X$  be separable. Then for each  $\varepsilon > 0$  there exist a measurable map  $\pi_{\varepsilon} : X \to M$  satisfying  $\rho(x, \pi_{\varepsilon}(x)) \leq \varepsilon + \rho(x, M)$  for all  $x \in M$ .

By condition C2),  $\mathscr{C}_{\delta}$  is separable with respect to the  $L^2$  norm, the continuous map  $\mathscr{C}_{\delta} \to L^2(\lambda)$  defined by  $b \to h = \mathbbm{1}_{b(\boldsymbol{x}^{(d)}) \leq x_d}$  implies that  $H_{\mathscr{C}_{\delta}}$  is separable. Then, for any  $B \in \mathcal{A}_n\left(L^2(\lambda)\right)$  and  $\varepsilon > 0$ , Lemma 16 implies that there exists a map  $\pi_{\varepsilon} : L^2(\lambda) \to H_{\mathscr{C}_{\delta}}$  such that

$$\|B(\boldsymbol{S}_h) - \pi_{\varepsilon}(B(\boldsymbol{S}_h))\|_{L^2(\mu_h)} \leq \varepsilon + \inf_{h \in H_{\mathscr{C}_{\varepsilon}}} \|B(\boldsymbol{S}_h) - h\|_{L^2(\mu_h)} \leq \varepsilon + \|B(\boldsymbol{S}_h) - h\|_{L^2(\mu_h)},$$

for all  $h \in H_{\mathscr{C}_{\delta}}$  and  $\mu_h$  as in C3). Therefore,

$$\begin{split} \widetilde{\mathcal{I}_{n}}(\mathscr{C}_{\delta}) &\leq \sup_{\substack{h \in H_{\mathscr{C}_{\delta}} \\ \mu_{h} \text{ has density as in (18)}}} \mathbb{E}_{\left\{\boldsymbol{x}_{i}\right\}_{i=1}^{n}} \overset{iid}{\sim} \mu_{h} } \left\| \pi_{\varepsilon}(B(\boldsymbol{S}_{h})) - h \right\|_{L^{2}(\mu_{h})}^{2} \\ &\leq \sup_{\substack{h \in H_{\mathscr{C}_{\delta}} \\ \mu_{h} \text{ satisfies C3)}}} \mathbb{E}_{\left\{\boldsymbol{x}_{i}\right\}_{i=1}^{n}} \overset{iid}{\sim} \mu_{h} } \left\| \pi_{\varepsilon}(B(\boldsymbol{S}_{h})) - h \right\|_{L^{2}(\mu_{h})}^{2} \\ &\leq \sup_{\substack{h \in H_{\mathscr{C}_{\delta}} \\ \mu_{h} \text{ satisfies C3)}}} \mathbb{E}_{\left\{\boldsymbol{x}_{i}\right\}_{i=1}^{n}} \overset{iid}{\sim} \mu_{h} } \left( \left\| \pi_{\varepsilon}(B(\boldsymbol{S}_{h})) - B(\boldsymbol{S}_{h}) \right\|_{L^{2}(\mu_{h})} + \left\| B(\boldsymbol{S}_{h}) - h \right\|_{L^{2}(\mu_{h})} \right)^{2} \\ &\leq \sup_{\substack{h \in H_{\mathscr{C}_{\delta}} \\ \mu_{h} \text{ satisfies C3)}}} \mathbb{E}_{\left\{\boldsymbol{x}_{i}\right\}_{i=1}^{n}} \overset{iid}{\sim} \mu_{h} } \left( \varepsilon + 2 \left\| B(\boldsymbol{S}_{h}) - h \right\|_{L^{2}(\mu_{h})} \right)^{2}, \end{split}$$

where  $\varepsilon \to 0$ , implies

$$\widetilde{\mathcal{I}_n}(\mathscr{C}_{\delta}) \leq 4 \sup_{\substack{h \in H_{\mathscr{C}_{\delta}} \\ \mu_h \text{ satisfies C3)}}} \mathbb{E}_{\{\boldsymbol{x}_i\}_{i=1}^n \overset{iid}{\sim} \mu_h} \left\| B(\boldsymbol{S}_h) - h \right\|_{L^2(\mu_h)}^2.$$

Furthermore, B was arbitrary and  $\mathscr{C}_{\delta} \subset \mathscr{C}$ , so

$$\widetilde{\mathcal{I}_n}(\mathscr{C}_{\delta}) \leq 4 \inf_{B \in \mathcal{A}_n(L^2(\lambda))} \sup_{\substack{h \in H_{\mathscr{C}_{\delta}} \\ \mu_h \text{ satisfies CS)}}} \mathbb{E}_{\{\boldsymbol{x}_i\}_{i=1}^n \overset{iid}{\sim} \mu_h} \left\| B(\boldsymbol{S}_h) - h \right\|_{L^2(\mu_h)}^2 \leq 4\mathcal{I}_n(\mathscr{C}),$$

and we obtain (37).

Since  $\mathscr{C}(\xi) \subseteq \mathscr{C}_{\delta}$ , it follows from (18) and the same argument used to get (37) with Lemma 16, that

$$\begin{split} 4\widetilde{\mathcal{I}_{n}}(\mathscr{C}_{\delta}) &= 4 \inf_{A \in \mathcal{A}_{n}(H_{\mathscr{C}_{\delta}})} \sup_{h_{b} \in H_{\mathscr{C}_{\delta}}} \mathbb{E}_{\left\{\boldsymbol{x}_{i}\right\}_{i=1}^{n}} \widetilde{\sim} \widetilde{\mu}_{h_{b}} \left\| A(\boldsymbol{S}_{h_{b}}) - h_{b} \right\|_{L^{2}(\widetilde{\mu}_{h_{b}})}^{2} \\ &\geq \inf_{A \in \mathcal{A}_{n}(H_{\mathscr{C}(\xi)})} \sup_{h_{b} \in H_{\mathscr{C}(\xi)}} \mathbb{E}_{\left\{\boldsymbol{x}_{i}\right\}_{i=1}^{n}} \widetilde{\sim} \widetilde{\mu}_{h_{b}} \left\| A(\boldsymbol{S}_{h_{b}}) - h_{b} \right\|_{L^{2}(\widetilde{\mu}_{h_{b}})}^{2}, \end{split}$$

where  $A \in \mathcal{A}_n(H_{\mathscr{C}(\xi)})$  implies  $A(\cdot) \in H_{\mathscr{C}(\xi)}$ , and there exists a unique  $b_A \in \mathscr{C}(\xi)$  such that  $A(\cdot) = h_{b_{A(\cdot)}}$ . Then, by Lemma 14, we get

$$\|A(S_{h_b}) - h_b\|_{L^2(\widetilde{\mu}_{h_b})}^2 = \frac{1}{\widetilde{\gamma} + 1} \|b_{A(S_{h_b})} - b\|_{L^{\widetilde{\gamma} + 1}([0,1]^{d-1})}^{\widetilde{\gamma} + 1} \ge \frac{1}{2\widetilde{\gamma}} \|b_{A(S_{h_b})} - b\|_{L^1([0,1]^{d-1})}^{\widetilde{\gamma} + 1}.$$

Therefore,

$$(8\widetilde{\gamma})\widetilde{\mathcal{I}_n}(\mathscr{C}_{\delta}) \ge \inf_{b_A \in \mathscr{C}(\xi)} \sup_{b \in \mathscr{C}(\xi)} \mathbb{E}_{\{\boldsymbol{x}_i\}_{i=1}^n \overset{iid}{\sim} \widetilde{\mu}_{h_b}} \left\| b_{A(\boldsymbol{S}_{h_b})} - b \right\|_{L^1([0,1]^{d-1})}^{\widetilde{\gamma}+1}. \tag{38}$$

By hypothesis,  $\mathscr{C}(\xi)$  is a  $\xi$ -packing set as in (P), with respect to the  $L^1$  norm, which means that

$$\|b_{A(\mathbf{S}_{h_b})} - b\|_{L^1([0,1]^{d-1})} > \xi, \quad \text{for all} \quad b_A, b \in \mathscr{C}(\xi).$$
 (39)

Then, by (38) and (39), together with  $\mathscr{C}^{\xi}=\left\{ \xi^{\widetilde{\gamma}}b:b\in\mathscr{C}(\xi)\right\}$ , we obtain

$$(8\widetilde{\gamma})\widetilde{\mathcal{I}}_{n}(\mathscr{C}_{\delta}) \geq \inf_{b_{A} \in \mathscr{C}(\xi)} \sup_{b \in \mathscr{C}(\xi)} \mathbb{E}_{\{\boldsymbol{x}_{i}\}_{i=1}^{n} \stackrel{iid}{\sim} \widetilde{\mu}_{h_{b}}} \left\| b_{A(\boldsymbol{S}_{h_{b}})} - b \right\|_{L^{1}([0,1]^{d-1})} \xi^{\widetilde{\gamma}}$$

$$= \inf_{b_{A} \in \mathscr{C}(\xi)} \sup_{b \in \mathscr{C}(\xi)} \mathbb{E}_{\{\boldsymbol{x}_{i}\}_{i=1}^{n} \stackrel{iid}{\sim} \widetilde{\mu}_{h_{b}}} \left\| \xi^{\widetilde{\gamma}} b_{A(\boldsymbol{S}_{h_{b}})} - \xi^{\widetilde{\gamma}} b \right\|_{L^{1}([0,1]^{d-1})}$$

$$= \inf_{b_{A} \in \mathscr{C}^{\xi}} \sup_{b \in \mathscr{C}^{\xi}} \mathbb{E}_{\{\boldsymbol{x}_{i}\}_{i=1}^{n} \stackrel{iid}{\sim} \widetilde{\mu}_{h_{b}}} \left\| b_{A(\boldsymbol{S}_{h_{b}})} - b \right\|_{L^{1}([0,1]^{d-1})}. \tag{40}$$

From (32),

$$2\widetilde{\gamma}^2 \left\| b_{A(\boldsymbol{S}_{h_b})} - b \right\|_{L^1([0,1]^{d-1})} \geq \rho_H^2(p_{h_{b_A(\boldsymbol{S}_{h_b})}}, p_{h_b}),$$

for all  $b_A, b \in \mathscr{C}^{\xi}$ , i.e. for all  $p_{h_{b_A(S_{h_b})}}, p_{h_b} \in \mathcal{P}_{\mathscr{C}^{\xi}}$ . We define the map  $E \in \mathcal{A}_n(\mathcal{P}_{\mathscr{C}^{\xi}})$  such that  $E(\cdot) := p_{A(\cdot)} = p_{h_{b_{A(\cdot)}}}$ , and (30), (40) imply

$$(16\widetilde{\gamma}^{3})\widetilde{\mathcal{I}_{n}}(\mathscr{C}_{\delta}) \geq \inf_{E \in \mathcal{A}_{n}(\mathcal{P}_{\mathscr{C}^{\xi}})} \sup_{p_{h_{b}} \in \mathcal{P}_{\mathscr{C}^{\xi}}} \mathbb{E}_{\mathbf{S}_{h_{b}}} \overset{iid}{\sim} p_{h_{b}}} \rho_{H}^{2}(E(\mathbf{S}_{h_{b}}), p_{h_{b}})$$

$$\geq \inf_{E \in \mathcal{A}_{n}(\mathcal{D}_{\mu})} \sup_{p' \in \mathcal{P}_{\mathscr{C}^{\xi}}} \mathbb{E}_{\mathbf{S}} \overset{iid}{\sim} p'} \rho_{H}^{2}(E(\mathbf{S}), p'), \tag{41}$$

since  $A_n(\mathcal{P}_{\mathscr{C}^{\xi}}) \subseteq A_n(\mathcal{P}_{\mathscr{C}_{\delta}}) \subseteq A_n(\mathcal{D}_{\mu})$ . In conclusion, (35) follows from (37) and (41).

Now we continue with the second part of the lemma, showing that (36) holds as follows. From Lemma 14 we have that

$$\frac{1}{2\widetilde{\gamma}} \|b_1 - b_2\|_{L^1([0,1]^{d-1})}^{\widetilde{\gamma}+1} \le \rho_H^2(p_{h_{b_1}}, p_{h_{b_2}}) \le 2\widetilde{\gamma}^2 \|b_1 - b_2\|_{L^1([0,1]^{d-1})}, \tag{42}$$

for all  $p_{h_{b_1}}, p_{h_{b_2}} \in \mathcal{P}_{\mathscr{C}_{\delta}}$ . Using the same argument as in (39) and the right hand side of (42), we obtain

$$\begin{split} \rho_H^2(p_{h_{b_1}}, p_{h_{b_2}}) &\leq 2\widetilde{\gamma}^2 \|b_1 - b_2\|_{L^1([0,1]^{d-1})} \\ &= 2\widetilde{\gamma}^2 \xi^{\widetilde{\gamma}} \|b_1' - b_2'\|_{L^1([0,1]^{d-1})} \\ &\leq 2\widetilde{\gamma}^2 \|b_1' - b_2'\|_{L^1([0,1]^{d-1})}^{\widetilde{\gamma}+1} \,, \end{split}$$

for all  $p_{h_{b_1}}, p_{h_{b_2}} \in \mathcal{P}_{\mathscr{C}^{\xi}}$  and all  $b_1', b_2' \in \mathscr{C}(\xi)$  associated with  $b_1, b_2 \in \mathscr{C}^{\xi}$ . Therefore,

$$\begin{split} M_{\mathcal{P}_{\mathscr{C}^{\xi}},\rho_{H}}(\varepsilon) &\leq M_{\mathscr{C}(\xi),L^{1}}\left(\varepsilon^{\frac{2}{\widetilde{\gamma}+1}}/(2\widetilde{\gamma}^{2})\right) \\ &\leq M_{\mathscr{C},L^{1}}\left(\varepsilon^{\frac{2}{\widetilde{\gamma}+1}}/(2\widetilde{\gamma}^{2})\right) \quad \text{for all} \quad \varepsilon > 0. \end{split} \tag{43}$$

Moreover, from the left hand side of (42),

$$\frac{1}{2\widetilde{\gamma}} \|b_1 - b_2\|_{L^1([0,1]^{d-1})}^{\widetilde{\gamma}+1} \le \rho_H^2(p_{h_{b_1}}, p_{h_{b_2}}),$$

for all  $p_{h_{b_1}}, p_{h_{b_2}} \in \mathcal{P}_{\mathscr{C}^{\xi}}$ , and

$$V_{\mathcal{P}_{\mathscr{C}^{\xi}},\rho_{H}}(\varepsilon) \geq V_{\mathscr{C}^{\xi},L^{1}}\left((2\widetilde{\gamma})\varepsilon^{\frac{2}{\widetilde{\gamma}+1}}\right) \geq V_{\mathscr{C},L^{1}}\left((2\widetilde{\gamma})\varepsilon^{\frac{2}{\widetilde{\gamma}+1}}\right) \quad \text{for all} \quad \varepsilon > 0. \tag{44}$$

Then, (7), (43) and (44) imply

$$V_{\mathscr{C},L^{1}}\left((2\widetilde{\gamma})\varepsilon^{\frac{2}{\widetilde{\gamma}+1}}\right) \leq V_{\mathcal{P}_{\mathscr{L}^{\xi}},\rho_{H}}(\varepsilon) \leq M_{\mathcal{P}_{\mathscr{L}^{\xi}},\rho_{H}}(\varepsilon) \leq M_{\mathscr{C},L^{1}}\left(\varepsilon^{\frac{2}{\widetilde{\gamma}+1}}/(2\widetilde{\gamma}^{2})\right),$$

for all  $\varepsilon > 0$ . Finally, from (7) and the above inequality, we conclude that

$$\begin{split} V_{\mathcal{P}_{\mathscr{C}^{\xi}},\rho_{H}}(\varepsilon) &\leq M_{\mathscr{C},L^{1}}\left(\varepsilon^{\frac{2}{\tilde{\gamma}+1}}/(2\tilde{\gamma}^{2})\right) \leq V_{\mathscr{C},L^{1}}\left(\varepsilon^{\frac{2}{\tilde{\gamma}+1}}/(4\tilde{\gamma}^{2})\right) \quad \text{an} \\ M_{\mathcal{P}_{\mathscr{C}^{\xi}},\rho_{H}}(\varepsilon) &\geq V_{\mathscr{C},L^{1}}\left((2\tilde{\gamma})\varepsilon^{\frac{2}{\tilde{\gamma}+1}}\right) \geq M_{\mathscr{C},L^{1}}\left((4\tilde{\gamma})\varepsilon^{\frac{2}{\tilde{\gamma}+1}}\right), \end{split}$$

for all  $\varepsilon > 0$ . So, (36) is fulfilled.

# B PROOF OF THE MAIN RESULTS

# B.1 Proof of Theorem 6

To begin with, we present (Petersen & Voigtlaender, 2021, Theorem 3.11), which is a simplified version of (Yang & Barron, 1999, Theorems 1 and 2).

**Lemma 17.** Let  $\emptyset \neq \mathcal{P} \subset D_{\mu}$ , and let  $V, M:(0,\infty) \to [0,\infty)$  be continuous from the right, non-increasing, and with the following properties:

1. 
$$V_{\mathcal{P},\rho_{KL}}(\varepsilon) \leq V(\varepsilon)$$
 for all  $\varepsilon > 0$ ;

2. 
$$(\varepsilon_n)_{n\in\mathbb{N}}\subset (0,\infty)$$
 is chosen such that  $\varepsilon_n^2=V(\varepsilon_n)/n$  for all  $n\in\mathbb{N}$ ;

3. 
$$M(\varepsilon) \leq M_{\mathcal{P},\rho_H}(\varepsilon)$$
 for all  $\varepsilon > 0$ ;

4. 
$$M(\varepsilon) > 2 \log 2$$
 for  $\varepsilon > 0$  small enough;

5. 
$$(\widetilde{\varepsilon}_n)_{n\in\mathbb{N}}\subset (0,\infty)$$
 is chosen such that  $M(\widetilde{\varepsilon}_n)=4n\varepsilon_n^2+2\log 2$  for all  $n\in N$ .

Then

$$\widehat{\varepsilon}_n^2/8 \le \inf_{E \in \mathcal{A}_n(\mathcal{D}_{\mu})} \sup_{p' \in \mathcal{P}} \mathbb{E}_{\mathbf{S}_{n'}^{iid}p'} \rho_H^2(E(\mathbf{S}), p') \le 2\varepsilon_n^2. \tag{45}$$

**Remark 18.** Let  $\emptyset \neq \mathcal{P} \subset D_{\mu}$ . In (Petersen & Voigtlaender, 2021, Page 15) it was shown that  $V_{\mathcal{P},\rho_{KL}}(\varepsilon) \leq V_{\mathcal{P},\rho_H}(\varepsilon/(2c))$  for all  $0 < \varepsilon < 1$  and a suitable absolute constant c > 1. Then, item 1 in Lemma 17 is fulfilled when  $V_{\mathcal{P},\rho_H}(\varepsilon/(2c)) \leq V(\varepsilon)$ .

Now, by hypothesis (10), Remark 18 and Lemma 15, it follows that

$$\begin{split} V_{\mathcal{P}_{\mathscr{C}^{\xi}},\rho_{KL}}(\varepsilon) &\leq V_{\mathcal{P}_{\mathscr{C}^{\xi}},\rho_{H}}(\varepsilon/(2c)) \leq V_{\mathscr{C},L^{1}}\left(\varepsilon^{\frac{2}{\widetilde{\gamma}+1}}/\left(4\widetilde{\gamma}^{2}(2c)^{\frac{2}{\widetilde{\gamma}+1}}\right)\right) \leq V(\varepsilon) \quad \text{and} \\ M_{\mathcal{P}_{\mathscr{C}^{\xi}},\rho_{H}}(\varepsilon) &\geq M_{\mathscr{C},L^{1}}\left((4\widetilde{\gamma})\varepsilon^{\frac{2}{\widetilde{\gamma}+1}}\right) \geq M(\varepsilon) \quad \text{for all} \quad \varepsilon > 0. \end{split} \tag{46}$$

Therefore, Lemma 17 holds with  $\mathcal{P} = \mathcal{P}_{\mathscr{C}^{\xi}}$ , since items 1 and 3 are satisfied with (46), and items 2, 4, and 5 are satisfied with (11). Then, using (35) and (45), we conclude that

$$\mathcal{I}_n(\mathscr{C}) \geq \frac{1}{(4\widetilde{\gamma})^3} \left[ \inf_{E \in \mathcal{A}_n(\mathcal{D}_{\boldsymbol{\mu}})} \sup_{p' \in \mathcal{P}_{\mathscr{C}^{\xi}}} \mathbb{E}_{\boldsymbol{S}} \underset{p'}{\overset{iid}{\sim}} p' \rho_H^2(E(\boldsymbol{S}), p') \right] \geq \left[ \frac{1}{8(4\widetilde{\gamma})^3} \right] \widetilde{\varepsilon}_n^2$$

and (12) is satisfied.

## B.2 PROOF OF COROLLARY 8

We consider each case as follows:

I) For all 
$$\varepsilon \in (0,1)$$
,  $\varepsilon^{-a} \lesssim M_{\mathscr{C},L^1}(\varepsilon) \lesssim \varepsilon^{-a}(1+\log(1/\varepsilon))$ . So, (7) implies

$$V_{\mathscr{C},L^1}(\varepsilon) \lesssim \varepsilon^{-a} (1 + \log(1/\varepsilon)) \quad \text{and} \quad \varepsilon^{-a} \lesssim M_{\mathscr{C},L^1}(\varepsilon).$$

To use Theorem 6, we set  $\widetilde{\varepsilon}_n, \varepsilon_n \in (0, 1)$ ,

$$\begin{split} V(\varepsilon) &\approx \left(\varepsilon^{\frac{2}{\overline{\gamma}+1}}\right)^{-a} \left(1 + \log\left(\varepsilon^{-\frac{2}{\overline{\gamma}+1}}\right)\right) \gtrsim V_{\mathscr{C},L^1}\left(\varepsilon^{\frac{2}{\overline{\gamma}+1}}\right) \quad \text{and} \\ M(\varepsilon) &\approx \left(\varepsilon^{\frac{2}{\overline{\gamma}+1}}\right)^{-a} \lesssim M_{\mathscr{C},L^1}\left(\varepsilon^{\frac{2}{\overline{\gamma}+1}}\right), \end{split}$$

such that (10) is satisfied. From (11),

$$n\varepsilon_n^2 \approx \left(\varepsilon_n^{\frac{2}{\tilde{\gamma}+1}}\right)^{-a} \left(1 + \log\left(\varepsilon_n^{-\frac{2}{\tilde{\gamma}+1}}\right)\right) \approx \left(\left(\widetilde{\varepsilon}_n^{\frac{2}{\tilde{\gamma}+1}}\right)^{-a} - 2\log 2\right)/4.$$
 (47)

Also,  $a \geq 1/2$  and  $\widetilde{\varepsilon}_n \leq 4\varepsilon_n$ . Then,

$$n \gtrsim 4n\varepsilon_n^2 + 2\log 2 \approx \left(\widetilde{\varepsilon}_n^{\frac{2}{\widetilde{\gamma}+1}}\right)^{-a} \gtrsim \varepsilon_n^{-\frac{1}{\widetilde{\gamma}+1}}.$$
 (48)

From (47) and (48),

$$n\varepsilon_n^2 \lesssim \left(\varepsilon_n^{\frac{2}{\widetilde{\gamma}+1}}\right)^{-a} \left(1+2\log n\right), \quad \text{i.e.} \quad \left(\varepsilon_n^2\right)^{1+\frac{a}{\widetilde{\gamma}+1}} \lesssim n^{-1} \left(1+2\log n\right).$$

Thus, 
$$\varepsilon_n^2 \lesssim \left[n^{-1} \left(1 + 2 \log n\right)\right]^{\frac{1}{1 + \left(\frac{a}{\gamma + 1}\right)}}$$
. By (47),

$$\left(\widehat{\varepsilon}_n^{\frac{2}{7+1}}\right)^{-a} \approx 4n\varepsilon_n^2 + 2\log 2 \lesssim n\left[n^{-1}\left(1+2\log n\right)\right]^{\frac{1}{1+\left(\frac{a}{7+1}\right)}} + 2\log 2$$

$$\lesssim n \left[ n^{-1} \left( 1 + 2 \log n \right) \right]^{\frac{1}{1 + \left( \frac{a}{\widetilde{\gamma} + 1} \right)}}.$$

Therefore,

$$\begin{split} \widetilde{\varepsilon}_n^2 \gtrsim n^{-\frac{1}{\left(\frac{a}{\gamma+1}\right)}} \left[n^{-1} \left(1 + 2\log n\right)\right]^{-\frac{1}{\left(1 + \left(\frac{a}{\gamma+1}\right)\right)\left(\frac{a}{\gamma+1}\right)}} \\ &= n^{-\frac{1}{1 + \left(\frac{a}{\gamma+1}\right)}} \left(1 + 2\log n\right)^{-\frac{1}{\left(1 + \left(\frac{a}{\gamma+1}\right)\right)\left(\frac{a}{\gamma+1}\right)}}, \end{split}$$

and (12) implies

$$\mathcal{I}_n(\mathscr{C}) \gtrsim \left[ n \left( 1 + 2 \log n \right)^{\frac{1}{\left(\frac{a}{\gamma+1}\right)}} \right]^{-\frac{1}{1+\left(\frac{a}{\gamma+1}\right)}}$$

II) We know that  $V_{\mathscr{C},L^1}(\varepsilon) \approx \varepsilon^{-a}$ , for all  $\varepsilon \in (0,1)$ . Then, to use Theorem 6, we choose

$$M(\varepsilon) \approx V(\varepsilon) \approx \left(\varepsilon^{\frac{2}{\widetilde{\gamma}+1}}\right)^{-a},$$

such that (10) is satisfied. So, (11) implies

$$\left(\varepsilon_n^2\right)^{1+\left(rac{a}{\overline{\gamma}+1}
ight)} pprox n^{-1}, \quad \text{i.e.} \quad \varepsilon_n^2 pprox n^{-\frac{1}{1+\left(rac{a}{\overline{\gamma}+1}
ight)}}.$$

Thus, using the above inequality and (11),

$$\left(\widetilde{\varepsilon}_{n}^{\frac{2}{\widetilde{\gamma}+1}}\right)^{-a} \approx 4n\varepsilon_{n}^{2} + 2\log 2 \lesssim n^{1-\frac{1}{1+\left(\frac{a}{\widetilde{\gamma}+1}\right)}} = n^{\frac{\left(\frac{a}{\widetilde{\gamma}+1}\right)}{1+\left(\frac{a}{\widetilde{\gamma}+1}\right)}}$$

and

$$\widetilde{\varepsilon}_n^2 \gtrsim n^{-\frac{1}{1+\left(\frac{a}{\widetilde{\gamma}+1}\right)}}.$$

In conclusion, from (12), it follows

$$\mathcal{I}_n(\mathscr{C}) \gtrsim n^{-\frac{1}{1+\left(\frac{a}{\overline{\gamma}+1}\right)}}.$$

# B.3 Proof of Corollary 9

We split the proof into the following cases:

I) It is well known that the space of Barron regular functions is Lipschitz, therefore  $\mathcal{B}_C$  satisfies C2) with  $\alpha=1$ . Furthermore, from Lemma 3, we have that item I) of Corollary 8 holds with a=2(d-1)/(d+1). Then,

$$\mathcal{I}_n(\mathcal{B}_C) \gtrsim \left[ n \left( 1 + 2 \log n \right)^{\frac{\max\{3,\gamma\}}{\left(\frac{2(d-1)}{d+1}\right)}} \right]^{-\frac{\max\{3,\gamma\}}{\max\{3,\gamma\} + \left(\frac{2(d-1)}{d+1}\right)}}.$$

II) By definition, the  $\alpha$ -Hölder continuous space satisfies C2). Then, with Lemma 4, we obtain that item II) of Corollary 8 is satisfied with  $a=(d-1)/\alpha$ . Therefore,

$$\mathcal{I}_n(\mathcal{H}_\alpha) \gtrsim n^{-\frac{\max\{3,\gamma/\alpha\}}{\max\{3,\gamma/\alpha\}+(d-1)/\alpha}}.$$

III) We know that  $C_B([0,1]^{d-1})$  is the set of all convex functions on  $[0,1]^{d-1}$  that are uniformly bounded by B. Therefore,  $C_B$  is Lipschitz and condition C2) is fulfilled with  $\alpha=1$ . Then, Lemma 5 implies that item II) of Corollary 8 holds with a=(d-1)/2. Thus,

$$\mathcal{I}_n(\mathcal{C}_B) \gtrsim n^{-\frac{2\max\{3,\gamma\}}{2\max\{3,\gamma\}+(d-1)}}.$$