
Boosted Histogram Transform for Regression (Supplementary Material)

Yuchao Cai^{*1,2} Hanyuan Hang^{*1} Hanfang Yang² Zhouchen Lin^{3,1}

A. Error Analysis

This section provides a more comprehensive error analysis for the theoretical results in Section 3. To be specific, we conduct approximation error analysis for the boosted regressor $f_{D,B}$ under the assumption that the Bayes decision function $f_{L,P}^*$ lies in the Höder spaces $C^{0,\alpha}$ and $C^{1,\alpha}$.

A.1. Error Analysis for $f_{L,P}^* \in C^{0,\alpha}$

First of all, we introduce some definitions and notations which will be used in the supplementary material. For a given histogram transform H , we write

$$f_{P,H} := \arg \min_{f \in \mathcal{F}_H} \mathcal{R}_{L,P}(f). \quad (\text{A.1})$$

In other words, $f_{P,H}$ is the function that minimizes the excess risk $\mathcal{R}_{L,P}(f)$ over the function set \mathcal{F}_H with the bin width $h \in [\underline{h}_0, \bar{h}_0]$. Then, elementary calculation yields

$$\begin{aligned} f_{P,H} &= \mathbb{E}_P(f_{L,P}^*(X) | A_H(x)) \\ &= \sum_{j \in \mathcal{I}_H} \frac{\int_{A_j} f_{L,P}^* dP_X}{P_X(A_j)} \cdot \mathbf{1}_{A_j} \\ &= \sum_{j \in \mathcal{I}_H} \frac{\int_{A_j} \mathbb{E}(Y|X) dP_X}{P_X(A_j)} \cdot \mathbf{1}_{A_j}. \end{aligned}$$

Moreover, we write

$$f_{D,H} := \arg \min_{f \in \mathcal{F}_H} \mathcal{R}_{L,D}(f) \quad (\text{A.2})$$

for the empirical version, which can be further presented as

$$f_{D,H} = \sum_{j \in \mathcal{I}_H} \frac{\sum_{i=1}^n Y_i \mathbf{1}_{A_j}(X_i)}{\sum_{i=1}^n \mathbf{1}_{A_j}(X_i)} \cdot \mathbf{1}_{A_j}.$$

^{*}Equal contribution ¹AI Lab, Samsung Research China - Beijing (SRC-B), Beijing, China ²School of Statistics, Renmin University of China, Beijing, China ³Key Lab. of Machine Perception (MoE), School of EECS, Peking University, Beijing, China. Correspondence to: Hanfang Yang <hyang@ruc.edu.cn>.

A.1.1. BOUNDING THE APPROXIMATION ERROR TERM

The following proposition shows that the L_2 distance between $f_{P,H}$ and $f_{L,P}^*$ behaves polynomial in the regularization parameter λ if we choose the bin width \underline{h}_0 appropriately.

Proposition 2 *Let the histogram transform H be defined as in (4) with bin width h satisfies Assumption 1. Furthermore, suppose that the Bayes decision function $f_{L,P}^* \in C^{0,\alpha}$. Then, for any fixed $\lambda > 0$, there holds*

$$\lambda h^{-2d} + \mathcal{R}_{L,P}(f_{P,H}) - \mathcal{R}_{L,P}^* \leq c \cdot \lambda^{\frac{\alpha}{\alpha+d}},$$

where c is some constant depending on α , d , and c_0 as in Assumption 1.

A.1.2. BOUNDING THE SAMPLE ERROR TERM

To derive bounds on the sample error of regularized empirical risk minimizers, let us briefly recall the definition of VC dimension measuring the complexity of the underlying function class.

Definition 3 (VC dimension) *Let \mathcal{B} be a class of subsets of \mathcal{X} and $A \subset \mathcal{X}$ be a finite set. The trace of \mathcal{B} on A is defined by $\{B \cap A : B \in \mathcal{B}\}$. Its cardinality is denoted by $\Delta^{\mathcal{B}}(A)$. We say that \mathcal{B} shatters A if $\Delta^{\mathcal{B}}(A) = 2^{\#(A)}$, that is, if for every $\tilde{A} \subset A$, there exists a $B \in \mathcal{B}$ such that $\tilde{A} = B \cap A$. For $k \in \mathbb{N}$, let*

$$m^{\mathcal{B}}(k) := \sup_{A \subset \mathcal{X}, \#(A)=k} \Delta^{\mathcal{B}}(A). \quad (\text{A.3})$$

Then, the set \mathcal{B} is a Vapnik-Chervonenkis class if there exists $k < \infty$ such that $m^{\mathcal{B}}(k) < 2^k$ and the minimal of such k is called the VC dimension of \mathcal{B} , and abbreviate as $\text{VC}(\mathcal{B})$.

To prove Lemma 4, we need the following fundamental lemma concerning with the VC dimension of purely random partitions, which follows the idea put forward by Breiman (2000) of the construction of purely random forest. To this end, let $p \in \mathbb{N}$ be fixed and π_p be a partition of \mathcal{X} with number of splits p and $\pi_{(p)}$ denote the collection of all partitions π_p .

Lemma 4 Let \mathcal{B}_p be defined by

$$\mathcal{B}_p := \left\{ B : B = \bigcup_{j \in J} A_j, J \subset \{0, 1, \dots, p\}, A_j \in \pi_p \right\}. \quad (\text{A.4})$$

Then the VC dimension of \mathcal{B}_p can be upper bounded by $dp + 2$.

To investigate the capacity property of continuous-valued functions, we need to introduce the concept *VC-subgraph class*. To this end, the *subgraph* of a function $f : \mathcal{X} \rightarrow \mathbb{R}$ is defined by

$$\text{sg}(f) := \{(x, t) : t < f(x)\}.$$

A class \mathcal{F} of functions on \mathcal{X} is said to be a *VC-subgraph class*, if the collection of all subgraphs of functions in \mathcal{F} , which is denoted by $\text{sg}(\mathcal{F}) := \{\text{sg}(f) : f \in \mathcal{F}\}$ is a VC class of sets in $\mathcal{X} \times \mathbb{R}$. Then the VC dimension of \mathcal{F} is defined by the VC dimension of the collection of the subgraphs, that is, $\text{VC}(\mathcal{F}) = \text{VC}(\text{sg}(\mathcal{F}))$.

Before we proceed, we also need to recall the definitions of the convex hull and VC-hull class. The symmetric *convex hull* $\text{Co}(\mathcal{F})$ of a class of functions \mathcal{F} is defined as the set of functions $\sum_{i=1}^m \alpha_i f_i$ with $\sum_{i=1}^m |\alpha_i| \leq 1$ and each f_i contained in \mathcal{F} . A set of measurable functions is called a *VC-hull class*, if it is in the pointwise sequential closure of the symmetric convex hull of a VC-class of functions.

We denote the function set \mathcal{F} as

$$\mathcal{F} := \bigcup_{H \sim \mathcal{P}_H} \mathcal{F}_H, \quad (\text{A.5})$$

which contains all the functions of \mathcal{F}_H induced by histogram transforms H with bin width \underline{h}_0 .

The following lemma presents the upper bound for the VC dimension of the function set \mathcal{F} .

Lemma 5 Let \mathcal{F} be the function set defined as in (A.5). Then \mathcal{F} is a VC-subgraph class with

$$\text{VC}(\mathcal{F}) \leq (d+1)2^{d+1}(\lfloor 2R\sqrt{d}/\underline{h}_0 \rfloor + 1)^d.$$

To further bound the capacity of the function sets, we need to introduce the following fundamental descriptions which enables an approximation of an infinite set by finite subsets.

Definition 6 (Covering Numbers) Let (\mathcal{X}, d) be a metric space, $A \subset \mathcal{X}$ and $\varepsilon > 0$. We call $A' \subset A$ an ε -net of A if for all $x \in A$ there exists an $x' \in A'$ such that $d(x, x') \leq \varepsilon$. Moreover, the ε -covering number of A is defined as

$$\mathcal{N}(A, d, \varepsilon) = \inf \left\{ n \geq 1 : \exists x_1, \dots, x_n \in \mathcal{X}, \right.$$

$$\left. \text{such that } A \subset \bigcup_{i=1}^n B_d(x_i, \varepsilon) \right\},$$

where $B_d(x, \varepsilon)$ denotes the closed ball in \mathcal{X} centered at x with radius ε .

The following lemma follows directly from Theorem 2.6.9 in Van der Vaart & Wellner (1996). For the sake of completeness, we present the proof in Section B.1.2.

Lemma 7 Let \mathbb{Q} be a probability measure on \mathcal{X} and

$$\mathcal{F} := \{f : \mathcal{X} \rightarrow \mathbb{R} : f \in [-M, M]\}.$$

Assume that for some fixed $\varepsilon > 0$ and $v > 0$, the covering number of \mathcal{F} satisfies

$$\mathcal{N}(\mathcal{F}, L_2(\mathbb{Q}), M\varepsilon) \leq c(1/\varepsilon)^v. \quad (\text{A.6})$$

Then there exists a universal constant c' such that

$$\log \mathcal{N}(\text{Co}(\mathcal{F}), L_2(\mathbb{Q}), M\varepsilon) \leq c' c^{2/(v+2)} \varepsilon^{-2v/(v+2)}.$$

The next theorem shows that covering numbers of \mathcal{F} grow at a polynomial rate.

Theorem 8 Let \mathcal{F} be a function set defined as in (A.5). Then there exists a universal constant $c < \infty$ such that for any $\varepsilon \in (0, 1)$ and any probability measure \mathbb{Q} , we have

$$\begin{aligned} \mathcal{N}(\mathcal{F}, L_2(\mathbb{Q}), M\varepsilon) \\ \leq c_0 (c_d / \underline{h}_0)^d \cdot (16e)^{(c_d / \underline{h}_0)^d} \varepsilon^{2(\underline{h}_0 / c_d)^d - 2}, \end{aligned}$$

where the constant $c_d := 2^{1+4/d} \cdot d^{1/2+1/d}$.

The following theorem gives an upper bound on the covering number of the VC-hull class $\text{Co}(\mathcal{F})$.

Theorem 9 Let \mathcal{F} be the function set defined as in (A.5). Then there exists a constant c_1 such that for any $\varepsilon \in (0, 1)$ and any probability measure \mathbb{Q} , there holds

$$\log \mathcal{N}(\text{Co}(\mathcal{F}), L_2(\mathbb{Q}), M\varepsilon) \leq c_1 \varepsilon^{2(\underline{h}_0 / c_d)^d - 2}. \quad (\text{A.7})$$

Next, let us recall the definition of entropy numbers.

Definition 10 (Entropy Numbers) Let (\mathcal{X}, d) be a metric space, $A \subset \mathcal{X}$ and $m \geq 1$ be an integer. The m -th entropy number of (A, d) is defined as

$$e_m(A, d) = \inf \left\{ \varepsilon > 0 : \exists x_1, \dots, x_{2^{m-1}} \in \mathcal{X} \right. \\ \left. \text{such that } A \subset \bigcup_{i=1}^{2^{m-1}} B_d(x_i, \varepsilon) \right\}.$$

Moreover, if (A, d) is a subspace of a normed space $(E, \|\cdot\|)$ and the metric d is given by $d(x, x') = \|x - x'\|$, $x, x' \in A$, we write $e_m(A, \|\cdot\|) := e_m(A, E) := e_m(A, d)$. Finally, if $S : E \rightarrow F$ is a bounded, linear operator between the normed space E and F , we denote $e_m(S) := e_m(SB_E, \|\cdot\|_F)$.

For a finite set $D \in \mathcal{X}^n$, we define the norm of an empirical L_2 -space by

$$\|f\|_{L_2(D)}^2 = \mathbb{E}_D |f|^2 := \frac{1}{n} \sum_{i=1}^n |f(x_i)|^2.$$

If E is the function space (9) and $D_X \in \mathcal{X}^n$, then the entropy number $e_m(\text{id} : E \rightarrow L_2(D_X))$ equals the m -th entropy number of the symmetric convex hull of the family $\{(f_i), f_i \in \mathcal{F}_i\}$, where $\text{id} : E \rightarrow L_2(D_X)$ denotes the identity map that assigns to every $f \in E$ the corresponding equivalence class in $L_2(D_X)$.

Now, we are able to present an oracle inequality for BHTR, which gives an upper bound for the sample error term.

Theorem 11 *Let the histogram transform H_n be defined as in (4) with bin width h_n satisfying Assumption 1. Furthermore, let $f_{D,B}$ be the BHTR regressor defined by (10) and $A(\lambda)$ be the corresponding approximation error defined by (11). Then for all $\tau > 0$, with probability $\mathbb{P}^n \otimes \mathbb{P}_H$ not less than $1 - 3e^{-\tau}$, we have*

$$\begin{aligned} & \Omega_\lambda(f) + \mathcal{R}_{L,D}(f_{D,B}) - \mathcal{R}_{L,P}^* \\ & \leq 12A(\lambda) + 3456M^2\tau/n + 3c'_0T^{2\delta'}\lambda_1^{-2\delta'}\lambda_2^{-1}n^{-2}, \end{aligned}$$

where c'_0 is a constant.

A.2. Error Analysis for $f_{L,P}^* \in C^{1,\alpha}$

A drawback to the analysis in $C^{0,\alpha}$ is that the usual Taylor expansion involved techniques for error estimation may not apply directly. As a result, we fail to prove the exact benefits of the boosting procedure. Therefore, in this subsection, we turn to the function space $C^{1,\alpha}$ consisting of smoother functions. To be specific, we study the convergence rates of $f_{D,B}$ to the Bayes decision function $f_{L,P}^* \in C^{1,\alpha}$. To this end, there is a point in introducing some notations.

For fixed $\underline{h}_0, \bar{h}_0 > 0$, let $\{H_t\}_{t=1}^T$ be histogram transforms with bin width $h_t \in [\underline{h}_0, \bar{h}_0]$, $t = 1, \dots, T$. Moreover, let $\{f_{P,H_t}\}_{t=1}^T$ and $\{f_{D,H_t}\}_{t=1}^T$ be defined as in (A.1) and (A.2), respectively. For $x \in \mathcal{X}$, we define

$$f_{P,E}(x) := \frac{1}{T} \sum_{t=1}^T f_{P,H_t}(x) \quad (\text{A.8})$$

and

$$f_{D,E}(x) := \frac{1}{T} \sum_{t=1}^T f_{D,H_t}(x). \quad (\text{A.9})$$

Then we make the error decomposition

$$\begin{aligned} & \mathbb{E}_{\nu_n} (\mathcal{R}_{L,P}(f_{D,E}) - \mathcal{R}_{L,P}^*) \\ & = \mathbb{E}_{\nu_n} \mathbb{E}_{P_X} (f_{D,E}(X) - f_{L,P}^*(X))^2 \\ & = \mathbb{E}_{\nu_n} \mathbb{E}_{P_X} (f_{D,E}(X) - f_{P,E}(X))^2 \\ & \quad + \mathbb{E}_{\nu_n} \mathbb{E}_{P_X} (f_{P,E}(X) - f_{L,P}^*(X))^2, \end{aligned} \quad (\text{A.10})$$

where $\nu_n := \mathbb{P}^n \otimes \mathbb{P}_H$. In particular, in the case that $T = 1$, i.e., for the base regressor HTR, we are concerned with the lower bound for $f_{D,H}$. We make the error decomposition

$$\begin{aligned} & \mathbb{E}_{\nu_n} (\mathcal{R}_{L,P}(f_{D,H}) - \mathcal{R}_{L,P}^*) \\ & = \mathbb{E}_{\nu_n} \mathbb{E}_{P_X} (f_{D,H}(X) - f_{L,P}^*(X))^2 \\ & = \mathbb{E}_{\nu_n} \mathbb{E}_{P_X} (f_{D,H}(X) - f_{P,H}(X))^2 \\ & \quad + \mathbb{E}_{\nu_n} \mathbb{E}_{P_X} (f_{P,H}(X) - f_{L,P}^*(X))^2. \end{aligned} \quad (\text{A.11})$$

It is important to note that both of the two terms on the right-hand side of (A.10) and (A.11) are data- and partition-independent due to the expectation with respect to D and H . Loosely speaking, the first error term corresponds to the expected estimation error of the estimators $f_{D,E}$ or $f_{D,H}$, while the second one demonstrates the expected approximation error.

A.2.1. UPPER BOUND FOR CONVERGENCE RATE OF BHTR

The following Lemma presents the explicit representation of $A_H(x)$ which will be used later in the proofs of Proposition 13.

Lemma 12 *Let the histogram transform H be defined as in (4) and A_H^1, A_H be as in (6) and (5), respectively. Then for any $x \in \mathbb{R}^d$, the set $A_H(x)$ can be represented as*

$$A_H(x) = \{x + (R \cdot S)^{-1}z : z \in [-b', 1 - b']\},$$

where $b' \sim \text{Unif}(0, 1)^d$.

The next proposition presents the upper bound of the L_2 distance between the ensemble regressor $f_{P,E}$ and the Bayes decision function $f_{L,P}^*$ in the Hölder space $C^{1,\alpha}$.

Proposition 13 *Let the histogram transform H be defined as in (4) with bin width h satisfying Assumption 1 and T be the number of iterations. Furthermore, let \mathbb{P}_X be the uniform distribution and $L_{\bar{h}_0}^-(x, y, t)$ be the restricted least squares loss defined as in (13). Moreover, let the Bayes decision function satisfy $f_{L,P}^* \in C^{1,\alpha}$. Then there holds*

$$\mathcal{R}_{L_{\bar{h}_0},P}(f_{P,E}) - \mathcal{R}_{L_{\bar{h}_0},P}^* \leq c_L^2 \bar{h}_0^{2(1+\alpha)} + \frac{1}{T} \cdot dc_L^2 \bar{h}_0^2 \quad (\text{A.12})$$

in expectation with respect to \mathbb{P}_H .

A.2.2. LOWER BOUND OF CONVERGENCE RATE OF HTR

The following two propositions present the lower bound of approximation error and sample error of HTR respectively.

Proposition 14 *Let the histogram transform H be defined as in (4) with bin width h satisfying Assumption 1 and $\bar{h}_0 \leq 1$. Moreover, let the regression model defined by*

$$Y := f(X) + \varepsilon \quad (\text{A.13})$$

with $f \in C^{1,\alpha}$. For a fixed constant $c_f \in (0, \infty)$, let \mathcal{A}_f be defined as

$$\mathcal{A}_f := \{x \in \mathbb{R}^d : \|\nabla f\|_\infty \geq c_f\} \quad (\text{A.14})$$

and N_1 be defined as

$$N_1 := \min \left\{ n \in \mathbb{N} : \bar{h}_{0,n} \leq \frac{R}{4\sqrt{d}} \right\}. \quad (\text{A.15})$$

Then for all $n > N_1$, there holds

$$\mathcal{R}_{L,P}(f_{P,H}) - \mathcal{R}_{L,P}^* \geq \frac{d}{12} \left(\frac{R}{2} \right)^d c_0^2 \mathbb{P}_X(\mathcal{A}_f) c_f^2 \cdot \bar{h}_0^2$$

in expectation with respect to P_H .

Proposition 15 *Let the histogram transform H be defined as in (4) with bin width h satisfying Assumption 1. Let the regression model be defined as in (A.13) with $f \in C^{1,\alpha}$. Moreover, assume that ε is independent of X such that $\mathbb{E}(\varepsilon|X) = 0$ and $\text{Var}(\varepsilon|X) =: \sigma^2 \leq 4M^2$ hold almost surely for some $M > 0$. Then there holds*

$$\begin{aligned} \mathcal{R}_{L,P}(f_{D,H}) - \mathcal{R}_{L,P}(f_{P,H}) \\ \geq 4R^d \sigma^2 (1 - 2e^{-1}) c_0^d \cdot \bar{h}_0^{-d} \cdot n^{-1} \end{aligned}$$

in expectation with respect to P^n , where the constant c_0 is as in Assumption 1.

B. Proofs

It is well-known that entropy numbers are closely related to the covering numbers. To be specific, entropy and covering numbers are in some sense inverse to each other. More precisely, for all constants $a > 0$ and $q > 0$, the implication

$$\begin{aligned} e_i(T, d) \leq ai^{-1/q}, \quad \forall i \geq 1 \\ \implies \ln \mathcal{N}(T, d, \varepsilon) \leq \ln(4)(a/\varepsilon)^q, \quad \forall \varepsilon > 0 \end{aligned} \quad (\text{B.1})$$

holds by Lemma 6.21 in Steinwart & Christmann (2008). Additionally, Exercise 6.8 in Steinwart & Christmann (2008) yields the opposite implication, namely

$$\begin{aligned} \ln \mathcal{N}(T, d, \varepsilon) < (a/\varepsilon)^q, \quad \forall \varepsilon > 0 \\ \implies e_i(T, d) \leq 3^{1/q} ai^{-1/q}, \quad \forall i \geq 1. \end{aligned} \quad (\text{B.2})$$

B.1. Proof for $f_{L,P}^* \in C^{0,\alpha}$

B.1.1. PROOF RELATED TO SECTION A.1.1

Proof [of Proposition 2] The assumption $f_{L,P}^* \in C^{0,\alpha}$ implies

$$\begin{aligned} \mathcal{R}_{L,P}(f_{P,H}) - \mathcal{R}_{L,P}^* \\ &= \|f_{P,H} - f_{L,P}^*\|_{L_2(P_X)}^2 \\ &= \left\| \sum_{j \in \mathcal{I}_H} \frac{\mathbf{1}_{A_j}(x)}{\mathbb{P}_X(A_j)} \int_{A_j} f_{L,P}^*(x') - f_{L,P}^*(x) d\mathbb{P}_X(x') \right\|_2^2 \\ &\leq \left\| \sum_{j \in \mathcal{I}_H} \frac{\mathbf{1}_{A_j}(x)}{\mathbb{P}_X(A_j)} \int_{A_j} |f_{L,P}^*(x') - f_{L,P}^*(x)| d\mathbb{P}_X(x') \right\|_2^2 \\ &\leq \left\| \sum_{j \in \mathcal{I}_H} \frac{\mathbf{1}_{A_j}(x)}{\mathbb{P}_X(A_j)} \int_{A_j} \|x' - x\|^\alpha d\mathbb{P}_X(x') \right\|_{L_2(P_X)}^2 \\ &\leq \left\| \sum_{j \in \mathcal{I}_H} \frac{\mathbf{1}_{A_j}(x)}{\mathbb{P}_X(A_j)} (\sqrt{d} \cdot \bar{h}_0)^\alpha \mathbb{P}_X(A_j) \right\|_{L_2(P_X)}^2 \\ &\leq (\sqrt{d} \cdot \bar{h}_0)^{2\alpha} \\ &\leq d^\alpha c_0^{-2\alpha} \underline{h}_0^{2\alpha}, \end{aligned}$$

where the last inequality follows from Assumption 2. Consequently we obtain

$$\begin{aligned} \lambda h^{-2d} + \mathcal{R}_{L,P}(f_{P,H}) - \mathcal{R}_{L,P}^* &\leq \lambda \underline{h}_0^{-2d} + d^\alpha c_0^{-2\alpha} \underline{h}_0^{2\alpha} \\ &\leq c \lambda^{\frac{\alpha}{\alpha+d}}, \end{aligned}$$

where the constant $c := d^{-\alpha/(2\alpha+2d)} c_0^{-\alpha/(\alpha+d)}$. \blacksquare

B.1.2. PROOF RELATED TO SECTION A.1.2

Proof [of Lemma 4] This proof is conducted from the perspective of geometric constructions.

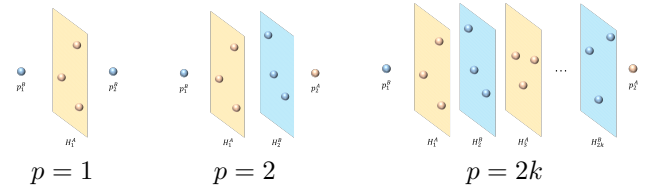


Figure 1. We take one case with $d = 3$ as an example to illustrate the geometric interpretation of the VC dimension. The yellow balls represent samples from class A, blue ones are from class B and slices denote the hyper-planes formed by samples.

We proceed by induction. Firstly, we concentrate on partition with the number of splits $p = 1$. Because of the dimension of the feature space is d , the smallest number of sample points that cannot be divided by $p = 1$ split is $d + 2$. Concretely, owing to the fact that d points can be used to

form $d - 1$ independent vectors and hence a hyperplane in a d -dimensional space, we might take the following case into consideration: There is a hyperplane consisting of d points all from one class, say class A , and two points p_1^B, p_2^B from the opposite class B located on the opposite sides of this hyperplane, respectively. We denote this hyperplane by H_1^A . In this case, points from two classes cannot be separated by one split (since the positions are p_1^B, H_1^A, p_2^B), so that we have $\text{VC}(\mathcal{B}_1) \leq d + 2$.

Next, when the partition is with the number of splits $p = 2$, we analyze in the similar way only by extending the above case a little bit. Now, we pick either of the two single sample points located on opposite side of the H_1^A , and add $d - 1$ more points from class B to it. Then, they together can form a hyperplane H_2^B parallel to H_1^A . After that, we place one more sample point from class A to the side of this newly constructed hyperplane H_2^B . In this case, the location of these two single points and two hyperplanes are $p_1^B, H_1^A, H_2^B, p_2^A$. Apparently, $p = 2$ splits cannot separate these $2d + 2$ points. As a result, we have $\text{VC}(\mathcal{B}_2) \leq 2d + 2$.

Inductively, the above analysis can be extended to the general case of number of splits $p \in \mathbb{N}$. In this manner, we need to add points continuously to form p mutually parallel hyperplanes where any two adjacent hyperplanes should be constructed from different classes. Without loss of generality, we consider the case for $p = 2k + 1$, $k \in \mathbb{N}$, where two points (denoted as p_1^B, p_2^B) from class B and $2k + 1$ alternately appearing hyperplanes form the space locations: $p_1^B, H_1^A, H_2^B, H_3^A, H_4^B, \dots, H_{(2k+1)}^A, p_2^B$. Accordingly, the smallest number of points that cannot be divided by p splits is $dp + 2$, leading to $\text{VC}(\mathcal{B}_p) \leq dp + 2$. This completes the proof. \blacksquare

Proof [of Lemma 5] Recall that for a histogram transform H , the set $\pi_H = (A_j)_{j \in \mathcal{I}_H}$ is a partition of B_R with the index set \mathcal{I}_H induced by H . The choice $k := \lfloor 2R\sqrt{d}/h_0 \rfloor + 1$ leads to the partition of B_R of the form $\pi_k := \{A_{i_1, \dots, i_d}\}_{i_j=1, \dots, k}$ with

$$\begin{aligned} A_{i_1, \dots, i_d} &:= \prod_{j=1}^d A_j \\ &:= \prod_{j=1}^d \left[-R + \frac{2R(i_j - 1)}{k}, -R + \frac{2Ri_j}{k} \right). \end{aligned} \quad (\text{B.3})$$

Obviously, we have $|A_{i_j}| \leq \frac{h_0}{\sqrt{d}}$. Let D be a data set of the form

$$D := \{(x_i, t_i) : x_i \in B_R, t_i \in [-M, M], i = 1, \dots, m\}$$

with

$$m := \#(D) = 2^{d+1}(d+1)(\lfloor 2R\sqrt{d}/h_0 \rfloor + 1)^d.$$

Then there exists at least one cell A with

$$\#(D \cap (A \times [-M, M])) \geq 2^{d+1}(d+1). \quad (\text{B.4})$$

Moreover, for any $x, x' \in A$, the construction of the partition (B.3) implies $\|x - x'\| \leq h_0$. Consequently, for any arbitrary histogram transform H and $A_j \in \pi_H$, at most one vertex of A_j lies in A , since the bin width of A_j is larger than h_0 . Therefore,

$$\begin{aligned} \Pi_{H|A} &:= \left\{ \bigcup_{j \in I} ((A_j \cap A) \times [-M, c_j]), I \subset \mathcal{I}_H \right\} \\ &\quad \bigcup \left\{ \bigcup_{j \in I} ((A_j \cap A) \times (c_j, M]), I \subset \mathcal{I}_H \right\} \end{aligned}$$

forms a partition of $A \times [-M, M]$ with $\#(\Pi_{H|A}) \leq 2^{d+1}$. It is easily seen that this partition can be generated by $2^{d+1} - 1$ splitting hyperplanes on the space $A \times [-M, M]$. In this way, Lemma 4 implies that $\Pi_{H|A}$ can only shatter a dataset with at most $(d+1)(2^{d+1} - 1) + 1$ elements. Thus (B.4) indicates that $\Pi_{H|A}$ fails to shatter $D \cap (A \times [-M, M])$. Therefore, the subgraphs of \mathcal{F}

$$\{ \{(x, t) : t < f(x)\}, f \in \mathcal{F} \}$$

cannot shatter the data set D as well. By Definition 3, we immediately get

$$\text{VC}(\mathcal{F}) \leq 2^{d+1}(d+1)(\lfloor 2R\sqrt{d}/h_0 \rfloor + 1)^d$$

and the assertion is thus proved. \blacksquare

Proof [of Lemma 7] Let \mathcal{F}_ε be an ε -net over \mathcal{F} . Then, for any $f \in \text{Co}(\mathcal{F})$, there exists an $f_\varepsilon \in \text{Co}(\mathcal{F}_\varepsilon)$ such that $\|f - f_\varepsilon\|_{L_2(\mathcal{Q})} \leq \varepsilon$. Therefore, we can assume without loss of generality that \mathcal{F} is finite.

Obviously, (A.6) holds for $1 \leq \varepsilon \leq c^{1/v}$. Let $v' := 1/2 + 1/v$ and $M' := c^{1/v}M$. Then (A.6) implies that for any $n \in \mathbb{N}$, there exists $f_1, \dots, f_n \in \mathcal{F}$ such that for any $f \in \mathcal{F}$, there exists an f_i such that

$$\|f - f_i\|_{L_2(\mathcal{Q})} \leq M'n^{-1/v}.$$

Therefore, for each $n \in \mathbb{N}$, we can find sets $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots \subset \mathcal{F}$ such that the set \mathcal{F}_n is a $M'n^{-1/v}$ -net over \mathcal{F} and $\#(\mathcal{F}_n) \leq n$.

In the following, we show by induction that for $q \geq 3 + v$ and $n, k \geq 1$, there holds

$$\log \mathcal{N}(\text{Co}(\mathcal{F}_{nk^q}), L_2(\mathcal{Q}), c_k M' n^{-v'}) \leq c'_k n, \quad (\text{B.5})$$

where c_k and c'_k are constants depending only on c and v such that $\sup_k \max\{c_k, c'_k\} < \infty$. The proof of (B.5) will be conducted by a nested induction argument.

Let us first consider the case $k = 1$. For a fixed n_0 , let $n \leq n_0$. Then for c_1 satisfying $c_1 M' n_0^{-v'}$ $\geq M$, there holds

$$\log \mathcal{N}(\text{Co}(\mathcal{F}_{nk^q}), L_2(\mathbb{Q}), c_k M' n^{-v'}) = 0,$$

which immediately implies (B.5). For a general $n \in \mathbb{N}$, let $m := n/\ell$ for large enough ℓ to be chosen later. Then for any $f \in \mathcal{F}_n \setminus \mathcal{F}_m$, there exists an $f^{(m)} \in \mathcal{F}_m$ such that

$$\|f - f^{(m)}\|_{L_2(\mathbb{Q})} \leq M' m^{-1/v}.$$

Let $\pi_m : \mathcal{F}_n \setminus \mathcal{F}_m \rightarrow \mathcal{F}_m$ be the projection operator. Then for any $f \in \mathcal{F}_n \setminus \mathcal{F}_m$, there holds

$$\|f - \pi_m f\|_{L_2(\mathbb{Q})} \leq M' m^{-1/v}$$

Therefore, for $\lambda_i, \mu_j \geq 0$ and $\sum_{i=1}^n \lambda_i = \sum_{j=1}^m \mu_j = 1$, we have

$$\sum_{i=1}^n \lambda_i f_i^{(n)} = \sum_{j=1}^m \mu_j f_j^{(m)} + \sum_{k=m+1}^n \lambda_k (f_k^{(n)} - \pi_m f_k^{(n)}).$$

Let \mathcal{G}_n be the set

$$\mathcal{G}_n := \{0\} \cup \{f - \pi_m f : f \in \mathcal{F}_n \setminus \mathcal{F}_m\}.$$

Then we have $\#(\mathcal{G}_n) \leq n$ and for any $g \in \mathcal{G}_n$, there holds

$$\|g\|_{L_2(\mathbb{Q})} \leq M' m^{-1/v}.$$

Moreover, we have

$$\text{Co}(\mathcal{F}_n) \subset \text{Co}(\mathcal{F}_m) + \text{Co}(\mathcal{G}_n). \quad (\text{B.6})$$

Applying Lemma 2.6.11 in Van der Vaart & Wellner (1996) with $\varepsilon := \frac{1}{2} c_1 m^{1/v} n^{-v'}$ to \mathcal{G}_n , we can find a $\frac{1}{2} c_1 M' n^{-v'}$ -net over $\text{Co}(\mathcal{G}_n)$ consisting of at most

$$(e + e n \varepsilon^2)^{2/\varepsilon^2} \leq \left(e + \frac{e c_1^2}{\ell^{2/v}} \right)^{8\ell^{2/v} c_1^{-2} n} \quad (\text{B.7})$$

elements.

Suppose that (B.5) holds for $k = 1$ and $n = m$. In other words, there exists a $c_1 M' m^{-v'}$ -net over $\text{Co}(\mathcal{F}_m)$ consisting of at most e^m elements, which partitions $\text{Co}(\mathcal{F}_m)$ into m -dimensional cells of diameter at most $2c_1 M' m^{-v'}$. Each of these cells can be isometrically identified with a subset of a ball of radius $c_1 M' m^{-v'}$ in \mathbb{R}^m and can be therefore further partitioned into

$$\left(\frac{3c_1 M' m^{-v'}}{\frac{1}{2} c_1 M' m^{-v'}} \right)^m = (6\ell^{v'})^{n/\ell}$$

cells of diameter $\frac{1}{2} c_1 M' n^{-v'}$. As a result, we get a $\frac{1}{2} c_1 M' n^{-v'}$ -net of $\text{Co}(\mathcal{F}_m)$ containing at most

$$e^m \cdot (6\ell^{v'})^{n/\ell} \quad (\text{B.8})$$

elements.

Now, (B.6) together with (B.7) and (B.8) yields that there exists a $c_1 M' n^{-v'}$ -net of $\text{Co}(\mathcal{F}_n)$ whose cardinality can be bounded by

$$e^{n/\ell} (6\ell^{v'})^{n/\ell} \left(e + \frac{e c_1^2}{\ell^{2/v}} \right)^{8\ell^{2/v} c_1^{-2} n} \leq e^n,$$

for suitable choices of c_1 and ℓ depending only on v . This concludes the proof of (B.5) for $k = 1$ and every $n \in \mathbb{N}$.

Let us consider a general $k \in \mathbb{N}$. Similarly as above, there holds

$$\text{Co}(\mathcal{F}_{nk^q}) \subset \text{Co}(\mathcal{F}_{n(k-1)^q}) + \text{Co}(\mathcal{G}_{n,k}), \quad (\text{B.9})$$

where the set $\mathcal{G}_{n,k}$ contains at most $n k^q$ elements with norm smaller than $M'(n(k-1)^q)^{-1/v}$. Applying Lemma 2.6.11 in Van der Vaart & Wellner (1996) to $\mathcal{G}_{n,k}$, we can find an $M' k^{-2} n^{-v'}$ -net over $\text{Co}(\mathcal{G}_{n,k})$ consisting of at most

$$(e + e k^{2q/v-4+q})^{2^{2q/v+1} k^{4-2q/v} n} \quad (\text{B.10})$$

elements. Moreover, by the induction hypothesis, we have a $c_{k-1} M' n^{-v'}$ -net over $\text{Co}(\mathcal{F}_{n(k-1)^q})$ consisting of at most

$$e^{c'_{k-1} n} \quad (\text{B.11})$$

elements. Using (B.9), (B.10), and (B.11), we obtain a $c_k M' n^{-v'}$ -net over $\text{Co}(\mathcal{F}_{nk^q})$ consisting of at most $e^{c'_k n}$ elements, where

$$c_k = c_{k-1} + \frac{1}{k^2},$$

$$c'_k = c'_{k-1} + 2^{2q/v+1} \frac{1 + \log(1 + k^{2q/v-4+q})}{k^{2q/v-4}}.$$

Form the elementary analysis we know that if $2q/v - 5 = 2$, then there exist constants c'_1, c'_2 , and c'_3 such that

$$\lim_{k \rightarrow \infty} c_k = c^{-1/v} n_0^{(v+2)/2v} + \sum_{i=2}^{\infty} 1/i^2 \leq c'_1 c^{-1/v} + c'_2,$$

$$\lim_{k \rightarrow \infty} c'_k = 1 + c \sum_{i=1}^{\infty} 2(2/i)^{2q/v} i^5 \leq c'_3.$$

Thus (B.5) is proved. Taking $\varepsilon := c_k M' n^{-v'}/M$ in (B.5), we get

$$\log \mathcal{N}(\text{Co}(\mathcal{F}_{nk^q}), L_2(\mathbb{Q}), M\varepsilon)$$

$$\leq c'_k c_k^{1/v'} (M')^{1/v'} M^{-1/v'} \varepsilon^{-1/v'}.$$

This together with

$$(M')^{1/v'} = (c^{1/v} M)^{1/v'} = c^{2/(v+2)} M^{1/v'}$$

yields

$$\log \mathcal{N}(\text{Co}(\mathcal{F}), L_2(\mathcal{Q}), M\varepsilon) \leq c' c^{2/(v+2)} \varepsilon^{-2v/(v+2)},$$

where the constant c' depends on the constants c_1'' , c_2'' and c_3'' . This finishes the proof. ■

Proof [of Theorem 8] We find the upper bound of $\text{VC}(\mathcal{F})$ satisfies

$$\begin{aligned} 2^{d+1}(d+1)(2R\sqrt{d}/\underline{h}_0 + 2)^d &\leq d \cdot 2^{d+2}(4R\sqrt{d}/\underline{h}_0)^d \\ &= (c_d R/\underline{h}_0)^d, \end{aligned}$$

where $c_d := 2^{1+4/d} \cdot d^{1/2+1/d}$. Then Theorem 2.6.7 in Van der Vaart & Wellner (1996) yields the assertion. ■

Proof [of Theorem 9] The assertion follows directly from Lemma 7 with

$$\begin{aligned} c &:= c_0(c_d/\underline{h}_0)^d \cdot (16e)^{(c_d/\underline{h}_0)^d}, \\ v &:= 2((c_d/\underline{h}_0)^d - 1). \end{aligned}$$

Let $\delta := (\underline{h}_0/c_d)^d$, then we have

$$\begin{aligned} c^{2/(v+2)} &= (c_0\delta^{-1}(16e)^{1/\delta})^\delta \\ &= 16e(c_0\delta^{-1})^\delta \\ &= 16e(c_0\delta^{-1})^\delta. \end{aligned}$$

Note that the function f defined by $f(\delta) := (c_0\delta^{-1})^\delta$ is continuous and

$$\lim_{\delta \rightarrow 0} f(\delta) = 1.$$

Then there exists a constant $M_d > 0$ such that $f(\delta) \leq M_d$ for all $0 < \delta \leq (1/c_d)^d$ if $\underline{h}_0 \leq 1$. Consequently, we have

$$\log \mathcal{N}(\text{Co}(\mathcal{F}), L_2(\mathcal{Q}), M\varepsilon) \leq 16ec' M_d \varepsilon^{2(\underline{h}_0 \cdot n/c_d)^d - 2}.$$

With $c_1 := 16ec' M_d$ we obtain the assertion. ■

Proof [of Theorem 11] Denote

$$r^* := \Omega_\lambda(f) + \mathcal{R}_{L,P}(f) - R_{L,P}^*,$$

and for $r > r^*$, we write

$$\begin{aligned} \mathcal{F}_r &:= \{f \in E : \Omega_\lambda(f) + \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* \leq r\}, \\ \mathcal{H}_r &:= \{L \circ f - L \circ f_{L,P}^* : f \in \mathcal{F}_r\}. \end{aligned}$$

Note that for $f \in \mathcal{F}_r$, we have $\lambda_1 \|f\|_E \leq r$, that is,

$$\sum_{i \in I} |w_i|^2 \leq r/\lambda_1.$$

Then, by the Cauchy-Schwarz inequality, we get

$$\sum_{i \in I} |w_i| \leq \left(T \sum_{i \in I} |w_i|^2 \right)^{1/2} \leq (rT/\lambda_1)^{1/2}.$$

Consequently, we have $\mathcal{F}_r \subset (rT/\lambda_1)^{1/2} B_E$. Since L is Lipschitz continuous with $|L|_1 \leq 4M$, we find

$$\begin{aligned} \mathbb{E}_{D \sim P^n} e_m(\mathcal{H}_r, L_2(\mathcal{D})) &\leq 4M \mathbb{E}_{D \sim P_X^n} e_m(\mathcal{F}_r, L_2(\mathcal{D})) \\ &\leq 8M (rT/\lambda_1)^{1/2} \mathbb{E}_{D \sim P_X^n} e_m(B_E, L_2(\mathcal{D})) \\ &\leq 8M (rT/\lambda_1)^{1/2} \mathbb{E}_{D \sim P_X^n} e_m(\text{Co}(\mathcal{F}), L_2(\mathcal{D})). \end{aligned}$$

Let $\delta := (\underline{h}_0/c_d)^d$, $\delta' := 1 - \delta$, and $a := c_1^{1/(2\delta')} M$. Then (A.7) together with (B.2) implies that

$$e_m(\text{Co}(\mathcal{F}), L_2(\mathcal{D})) \leq (3c_1)^{1/(2\delta')} M i^{-1/(2\delta')}$$

Taking expectation with respect to P^n , we get

$$\mathbb{E}_{D \sim P_X^n} e_m(\text{Co}(\mathcal{F}), L_2(\mathcal{D})) \leq c_2 i^{-1/(2\delta')}, \quad (\text{B.12})$$

where $c_2 := (3c_1)^{1/(2\delta')} M$. Moreover, we easily find

$$\lambda_2 h^{-2d} = \Omega(h) \leq \Omega_\lambda(f) + \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* \leq r,$$

which yields

$$\underline{h}_0^{-1} \leq (r/\lambda_2)^{1/(2d)}.$$

Therefore, if $\underline{h}_0 \leq 1$, then we have $r \geq \lambda_2 \geq 1$ and (B.12) can be further estimated by

$$\mathbb{E}_{D \sim P_X^n} e_m(\text{Co}(\mathcal{F}_H), L_2(\mathcal{D})) \leq c_2 (r/\lambda_2)^{1/(4\delta')} i^{-1/(2\delta')},$$

which leads to

$$\begin{aligned} \mathbb{E}_{D \sim P_X^n} e_m(\mathcal{H}_r, L_2(\mathcal{D})) &\leq 8c_2 M (rT/\lambda_1)^{1/2} (r/\lambda_2)^{1/(4\delta')} i^{-1/(2\delta')}. \end{aligned}$$

For the least square loss, the supremum bound

$$L(x, y, t) \leq 4M^2, \quad \forall (x, y) \in \mathcal{X} \times \mathcal{Y}, t \in [-M, M],$$

and the variance bound

$$\mathbb{E}(L \circ g - L \circ f_{L,P}^*)^2 \leq V(\mathbb{E}(L \circ g - L \circ f_{L,P}^*))^\vartheta$$

holds for $V = 16M^2$ and $\vartheta = 1$. Therefore, for $h \in \mathcal{H}_r$, we have

$$\|h\|_\infty \leq 8M^2, \quad \mathbb{E}_P h^2 \leq 16M^2 r.$$

Then Theorem 7.16 in Steinwart & Christmann (2008) with $a := 8c_2 M (rT/\lambda_1)^{1/2} (r/\lambda_2)^{1/(4\delta')}$ yields that there exist a constant $c'_0 > 0$ such that

$$\mathbb{E}_{D \sim P^n} \text{Rad}_D(\mathcal{H}_r, n)$$

$$\begin{aligned} &\leq c'_0 \max \left\{ r^{3/4} T^{\delta'/2} \lambda_1^{-\delta'/2} \lambda_2^{-1/4} n^{-1/2}, \right. \\ &\quad \left. r^{(2\delta'+1)/(2\delta'+2)} (T/\lambda_1)^{\delta'/(1+\delta')} \right. \\ &\quad \left. \cdot \lambda_2^{1/(2+2\delta')} n^{-1/(1+\delta')} \right\} \\ &=: \varphi_n(r). \end{aligned}$$

Simple algebra shows that the condition $\varphi_n(4r) \leq 2\sqrt{2}\varphi_n(r)$ is satisfied. Since $2\sqrt{2} < 4$, similar arguments show that there still hold the statements of the Peeling Theorem 7.7 in [Steinwart & Christmann \(2008\)](#). Consequently, Theorem 7.20 in [Steinwart & Christmann \(2008\)](#) can also be applied, if the assumptions on φ_n and r are modified to $\varphi_n(4r) \leq 2\sqrt{2}\varphi_n(r)$ and $r \geq \max\{75\varphi_n(r), 1152M^2\tau/n, r^*\}$, respectively. It is easy to verify that the condition $r \geq 75\varphi_n(r)$ is satisfied if

$$r \geq c'_0 T^{2\delta'} \lambda_1^{-2\delta'} \lambda_2^{-1} n^{-2},$$

where c'_0 is a constant, which yields the assertion. \blacksquare

B.1.3. PROOF RELATED TO SECTION 3.1

Proof [of Theorem 1] It is easy to see that $f_{P,E}$ defined by (A.8) satisfies $f_{P,E} \in E$ and $\lambda_1 \|f_{P,E}\|_E \leq \lambda_1/T$. Moreover, by Jensen's inequality and Proposition 2, we have

$$\begin{aligned} \mathcal{R}_{L,P}(f_{P,E}) - \mathcal{R}_{L,P}^* &= \int_{\mathcal{X}} \left(\frac{1}{T} \sum_{t=1}^T f_{P,H_t} - f_{L,P}^* \right)^2 dP_X \\ &\leq \frac{1}{T} \sum_{t=1}^T \int_{\mathcal{X}} (f_{P,H_t} - f_{L,P}^*)^2 dP_X \\ &= \frac{1}{T} \sum_{t=1}^T \mathcal{R}_{L,P}(f_{P,H_t}) - \mathcal{R}_{L,P}^* \\ &\leq d^\alpha c_0^{-2\alpha} \underline{L}_0^{2\alpha}. \end{aligned}$$

Consequently we get

$$\begin{aligned} A(\lambda) &= \inf_{f \in E} \lambda_1 \|f\|_E + \lambda_2 \Omega(h) + \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* \\ &\leq \lambda_1 \|f_{P,E}\|_E + \lambda_2 \Omega(h) + \mathcal{R}_{L,P}(f_{P,E}) - \mathcal{R}_{L,P}^* \\ &\leq \lambda_1/T + c\lambda_2^{\frac{\alpha}{\alpha+d}}. \end{aligned}$$

Then, Theorem 11 implies that with probability $\mathbb{P} \otimes \mathbb{P}_H$ not less than $1 - 3e^{-\tau}$, there holds

$$\begin{aligned} &\lambda_1 \|f\|_E + \lambda_2 \Omega(h) + \mathcal{R}_{L,D}(f_{D,B}) - \mathcal{R}_{L,P}^* \\ &\leq 6\lambda_1/T + 6c\lambda_2^{\frac{\alpha}{\alpha+d}} + 3c'_0 T^{2\delta'} \lambda_1^{-2\delta'} \lambda_2^{-1} n^{-2} \\ &\quad + 3456M^2\tau/n, \end{aligned} \tag{B.13}$$

where c and c'_0 are constants defined as in Proposition 2 and Theorem 11. Minimizing the right hand side of (B.13), we

get

$$\mathcal{R}_{L,P}(f_{D,B}) - \mathcal{R}_{L,P}^* \leq c'' n^{-\frac{2\alpha}{(4-2\delta)\alpha+d}},$$

if we choose

$$\begin{aligned} \lambda_{1,n} &:= n^{-\frac{2\alpha}{(4-2\delta)\alpha+d}}, \\ \lambda_{2,n} &:= n^{-\frac{2(\alpha+d)}{(4-2\delta)\alpha+d}}, \\ h_{0,n} &:= n^{-\frac{1}{(4-2\delta)\alpha+d}}, \end{aligned}$$

where c'' is a constant depending on c, c'_0, d, M, R and T . Thus, the assertion is proved. \blacksquare

B.2. Proof for $f_{L,P}^* \in C^{1,\alpha}$

B.2.1. PROOF RELATED TO SECTION A.2.1

Proof [of Lemma 12] For any $x \in \mathbb{R}^d$, we define $b' := H(x) - \lfloor H(x) \rfloor \in \mathbb{R}^d$. Then we have $b' \sim \text{Unif}(0, 1)^d$ according to the definition of H . For any $x' \in A'_H(x)$, we define

$$z := H(x') - H(x) = (R \cdot S)(x' - x).$$

Then we have

$$x' = x + (R \cdot S)^{-1}z.$$

Moreover, since

$$\lfloor H(x') \rfloor = \lfloor H(x) \rfloor,$$

we have $z \in [-b', 1 - b']$. \blacksquare

Proof [of Proposition 13] According to the generation process, the histogram transforms $\{H_t\}_{t=1}^T$ are i.i.d. Therefore, for any $x \in B_R$, the expected approximation error term can be decomposed as follows:

$$\begin{aligned} &\mathbb{E}_{P_H} (f_{P,E}(x) - f_{L,P}^*(x))^2 \\ &= \mathbb{E}_{P_H} \left((f_{P,E}(x) - \mathbb{E}_{P_H}(f_{P,E}(x))) \right. \\ &\quad \left. + (\mathbb{E}_{P_H}(f_{P,E}(x)) - f_{L,P}^*(x)) \right)^2 \\ &= \text{Var}(f_{P,E}(x)) + (\mathbb{E}_{P_H}(f_{P,E}(x)) - f_{L,P}^*(x))^2 \\ &= \frac{1}{T} \cdot \text{Var}_{P_H}(f_{P,H_1}(x)) + (\mathbb{E}_{P_H}(f_{P,H_1}(x)) - f_{L,P}^*(x))^2. \end{aligned} \tag{B.14}$$

In the following, for the simplicity of notations, we drop the subscript of H_1 and write H instead of H_1 when there is no confusion.

For the first term in (B.14), the assumption $f_{L,P}^* \in C^{1,\alpha}$ implies

$$\text{Var}_{P_H}(f_{P,H}(x))$$

$$\begin{aligned}
 &= \mathbb{E}_{P_H} (f_{P,H}(x) - \mathbb{E}_{P_H}(f_{P,H}(x)))^2 \\
 &\leq \mathbb{E}_{P_H} (f_{P,H}(x) - f_{L,P}^*(x))^2 \\
 &= \mathbb{E}_{P_H} \left(\frac{1}{\mu(A_H(x))} \int_{A_H(x)} f_{L,P}^*(x') dx' - f_{L,P}^*(x) \right)^2 \\
 &= \mathbb{E}_{P_H} \left(\frac{1}{\mu(A_H(x))} \int_{A_H(x)} (f_{L,P}^*(x') - f_{L,P}^*(x)) dx' \right)^2 \\
 &\leq \mathbb{E}_{P_H} (c_L \text{diam}(A_H(x)))^2 \\
 &\leq c_L^2 d \bar{h}_0^2. \tag{B.15}
 \end{aligned}$$

We now consider the second term in (B.14). Lemma 12 implies that for any $x' \in A_H(x)$, there exist a random vector $u \sim \text{Unif}[0, 1]^d$ and a vector $v \in [0, 1]^d$ such that

$$x' = x + S^{-1}R^\top(-u + v). \tag{B.16}$$

Therefore, we have

$$\begin{aligned}
 dx' &= \det \left(\frac{dx'}{dv} \right) dv \\
 &= \det \left(\frac{d(x + S^{-1}R^\top(-u + v))}{dv} \right) dv \\
 &= \det(RS^{-1}) dv \\
 &= \left(\prod_{i=1}^d h_i \right) dv. \tag{B.17}
 \end{aligned}$$

Taking the first-order Taylor expansion of $f_{L,P}^*(x')$ at x , we get

$$\begin{aligned}
 &f_{L,P}^*(x') - f_{L,P}^*(x) \\
 &= \int_0^1 (\nabla f_{L,P}^*(x + t(x' - x)))^\top (x' - x) dt. \tag{B.18}
 \end{aligned}$$

Moreover, we obviously have

$$\nabla f_{L,P}^*(x)^\top (x' - x) = \int_0^1 \nabla f_{L,P}^*(x)^\top (x' - x) dt. \tag{B.19}$$

Thus, (B.18) and (B.19) imply that for any $f_{L,P}^* \in C^{1,\alpha}$, there holds

$$\begin{aligned}
 &|f_{L,P}^*(x') - f_{L,P}^*(x) - \nabla f_{L,P}^*(x)^\top (x' - x)| \\
 &= \left| \int_0^1 (\nabla f_{L,P}^*(x + t(x' - x)) - \nabla f_{L,P}^*(x))^\top (x' - x) dt \right| \\
 &\leq \int_0^1 c_L (t\|x' - x\|_2)^\alpha \|x' - x\|_2 dt \\
 &\leq c_L \|x' - x\|_2^{1+\alpha}.
 \end{aligned}$$

This together with (B.16) yields

$$|f_{L,P}^*(x') - f_{L,P}^*(x) - \nabla f_{L,P}^*(x)^\top S^{-1}R^\top(-u + v)|$$

$$\leq c_L \bar{h}_0^{1+\alpha}$$

and consequently there exists a constant $c_\alpha \in [-c_L, c_L]$ such that

$$\begin{aligned}
 &f_{L,P}^*(x') - f_{L,P}^*(x) \\
 &= \nabla f_{L,P}^*(x)^\top S^{-1}R^\top(-u + v) + c_\alpha \bar{h}_0^{1+\alpha}. \tag{B.20}
 \end{aligned}$$

Therefore, there holds

$$\begin{aligned}
 f_{P,H}(x) &= \frac{1}{P_X(A_H(x))} \int_{A_H(x)} f_{L,P}^*(x') dx' \\
 &= \frac{1}{\mu(A_H(x))} \int_{A_H(x)} f_{L,P}^*(x') dx'.
 \end{aligned}$$

This together with (B.20) and (B.17) yields

$$\begin{aligned}
 &f_{P,H}(x) - f_{L,P}^*(x) \\
 &= \frac{1}{\mu(A_H(x))} \int_{A_H(x)} f_{L,P}^*(x') dx' - f_{L,P}^*(x) \\
 &= \frac{1}{\mu(A_H(x))} \int_{A_H(x)} (f_{L,P}^*(x') - f_{L,P}^*(x)) dx' \\
 &= \frac{\prod_{i=1}^d h_i}{\mu(A_H(x))} \int_{[0,1]^d} \left(\nabla f_{L,P}^*(x)^\top S^{-1}R^\top(-u + v) \right. \\
 &\quad \left. + c_\alpha \bar{h}_0^{1+\alpha} \right) dv \\
 &= \left(\int_{[0,1]^d} (-u + v)^\top dv \right) RS^{-1} \nabla f_{L,P}^*(x) + c_\alpha \bar{h}_0^{1+\alpha} \\
 &= \left(\frac{1}{2} - u \right)^\top RS^{-1} \nabla f_{L,P}^*(x) + c_\alpha \bar{h}_0^{1+\alpha}. \tag{B.21}
 \end{aligned}$$

Since the random variables $(u_i)_{i=1}^d$ are independent and identically distributed as $\text{Unif}[0, 1]$, we have

$$\mathbb{E}_{P_H} \left(\frac{1}{2} - u_i \right) = 0, \quad i = 1, \dots, d. \tag{B.22}$$

Combining (B.21) with (B.22), we obtain

$$\mathbb{E}_{P_H} (f_{P,H}(x) - f_{L,P}^*(x)) = c_\alpha \bar{h}_0^{1+\alpha} \tag{B.23}$$

and consequently

$$(\mathbb{E}_{P_H}(f_{P,H}(x)) - f_{L,P}^*(x))^2 \leq c_L^2 \bar{h}_0^{2(1+\alpha)}. \tag{B.24}$$

Combining (B.14) with (B.24) and (B.15), we obtain

$$\mathbb{E}_{P_H} (f_{P,E}(x) - f_{L,P}^*(x))^2 \leq c_L^2 \bar{h}_0^{2(1+\alpha)} + \frac{1}{T} \cdot dc_L^2 \bar{h}_0^2.$$

Taking expectation with respect to P_X , we get

$$\mathbb{E}_{P_H} (\mathcal{R}_{L,P}(f_{P,E}) - \mathcal{R}_{L,P}^*) \leq c_L^2 \bar{h}_0^{2(1+\alpha)} + \frac{1}{T} \cdot dc_L^2 \bar{h}_0^2,$$

which completes the proof. \blacksquare

B.2.2. LOWER BOUND OF APPROXIMATION ERROR FOR HTR

Proof [of Proposition 14] Recall that the regression model is defined as $Y = f(X) + \varepsilon$. Considering the case when X follows the uniform distribution, for any $x = (x_1, \dots, x_d) \in \mathcal{X}$, we have

$$\begin{aligned} f_{\mathbb{P}, H}(x) &= \frac{1}{\mathbb{P}_X(A_H(x))} \int_{A_H(x)} f(x') dx' \\ &= \frac{1}{\mu(A_H(x))} \int_{A_H(x)} f(x') dx'. \end{aligned}$$

Then we get

$$\begin{aligned} &(f_{\mathbb{P}, H}(x) - f(x))^2 \\ &= \left(f(x) - \frac{1}{\mu(A_H(x))} \int_{A_H(x)} f(x') dx' \right)^2 \\ &= \frac{1}{\mu(A_H(x))^2} \left(\int_{A_H(x)} f(x') - f(x) dx' \right)^2. \end{aligned}$$

Lemma 12 implies that for any $x' \in A_H(x)$, there exists a random vector $u \sim \text{Unif}[0, 1]^d$ and a vector $v \in [0, 1]^d$ such that

$$x' = x + S^{-1}R^\top(-u + v). \quad (\text{B.25})$$

Therefore, we have

$$\begin{aligned} dx' &= \det\left(\frac{dx'}{dv}\right) dv \\ &= \det\left(\frac{d(x + S^{-1}R^\top(-u + v))}{dv}\right) dv \\ &= \det(RS^{-1}) dv \\ &= \left(\prod_{i=1}^d h_i\right) dv. \end{aligned}$$

Moreover, (B.20) yields that there exists a constant $c_\alpha \in [-c_L, c_L]$ such that

$$f(x') - f(x) = \nabla f(x)^\top S^{-1}R^\top(-u + v) + c_\alpha \bar{h}_0^{-1+\alpha}.$$

Taking expectation with regard to \mathbb{P}_H and \mathbb{P}_X , we get

$$\begin{aligned} &\mathbb{E}_{\mathbb{P}_X} (f_{\mathbb{P}, H}(X) - f(X))^2 \\ &\geq \mathbb{E}_{\mathbb{P}_X} (f_{\mathbb{P}, H}(X) - f_{L, \mathbb{P}}^*(X))^2 \mathbf{1}_{B_{R, \sqrt{d} \cdot \bar{h}_0}^+}(X) \\ &= \int_{B_{R, \sqrt{d} \cdot \bar{h}_0}^+} (f_{\mathbb{P}, H}(x) - f_{L, \mathbb{P}}^*(x))^2 d\mathbb{P}_X \\ &= \int_{B_{R, \sqrt{d} \cdot \bar{h}_0}^+} \frac{1}{\mu(A_H(x))^2} \left(\int_{A_H(x)} \nabla f(x)^\top S^{-1}R^\top(-u + v) + c_\alpha \bar{h}_0^{-1+\alpha} dy \right)^2 d\mathbb{P}_X \end{aligned}$$

$$\begin{aligned} &= \int_{B_{R, \sqrt{d} \cdot \bar{h}_0}^+} \frac{(\prod_{i=1}^d h_i)^2}{\mu(A_H(x))^2} \left(\int_{[0, 1]^d} (-u + v)^\top dv \right. \\ &\quad \left. \cdot RS^{-1}\nabla f(x) + c_\alpha \bar{h}_0^{-1+\alpha} \right)^2 d\mathbb{P}_X \\ &= \int_{B_{R, \sqrt{d} \cdot \bar{h}_0}^+} \left(\left(\frac{1}{2} - u \right)^\top RS^{-1}\nabla f(x) + c_\alpha \bar{h}_0^{-1+\alpha} \right)^2 d\mathbb{P}_X \\ &= \int_{B_{R, \sqrt{d} \cdot \bar{h}_0}^+} \left(\sum_{i=1}^d \left(\frac{1}{2} - u_i \right) \sum_{j=1}^d R_{ij} h_j \frac{\partial f}{\partial x_j} \right. \\ &\quad \left. + c_\alpha \bar{h}_0^{-1+\alpha} \right)^2 d\mathbb{P}_X. \quad (\text{B.26}) \end{aligned}$$

Since the random variables $(u_i)_{i=1}^d$ are i.i.d. as $\text{Unif}[0, 1]$, we have

$$\mathbb{E}_{\mathbb{P}_H} \left(\frac{1}{2} - u_i \right) = 0, \quad i = 1, \dots, d, \quad (\text{B.27})$$

and

$$\mathbb{E}_{\mathbb{P}_H} \left(\frac{1}{2} - u_i \right)^2 = \frac{1}{12}, \quad i = 1, \dots, d. \quad (\text{B.28})$$

Consequently, we have

$$\begin{aligned} &\mathbb{E}_{\mathbb{P}_H} \mathbb{E}_{\mathbb{P}_X} (f_{\mathbb{P}, H}(X) - f(X))^2 \\ &= \int_{B_{R, \sqrt{d} \cdot \bar{h}_0}^+} \mathbb{E}_{\mathbb{P}_H} \sum_{i=1}^d \left(\frac{1}{2} - u_i \right)^2 \left(\sum_{j=1}^d R_{ij} h_j \frac{\partial f}{\partial x_j} \right)^2 d\mathbb{P}_X. \end{aligned}$$

Moreover, the orthogonality (1) of the rotation matrix R tells us that

$$\sum_{i=1}^d R_{ij} R_{ik} = \begin{cases} 1, & \text{if } j = k, \\ 0, & \text{if } j \neq k \end{cases} \quad (\text{B.29})$$

and consequently we have

$$\begin{aligned} &\sum_{i=1}^d \sum_{j \neq k} R_{ij} R_{ik} h_j h_k \cdot \frac{\partial f(x)}{\partial x_j} \cdot \frac{\partial f(x)}{\partial x_k} \\ &= \sum_{j \neq k} h_j h_k \cdot \frac{\partial f(x)}{\partial x_j} \cdot \frac{\partial f(x)}{\partial x_k} \sum_{i=1}^d R_{ij} R_{ik} = 0. \quad (\text{B.30}) \end{aligned}$$

For any $n > N_1$ with N_1 as in (A.15), we have

$$(R - 2\sqrt{d} \cdot \bar{h}_0)^d \geq (R/2)^d.$$

Consequently, (B.29) and (B.30) imply that

$$\int_{B_{R, \sqrt{d} \cdot \bar{h}_0}^+} \mathbb{E}_{\mathbb{P}_H} \sum_{i=1}^d \left(\frac{1}{2} - u_i \right)^2 \left(\sum_{j=1}^d R_{ij} h_j \frac{\partial f}{\partial x_j} \right)^2 d\mathbb{P}_X$$

$$\begin{aligned}
 &= \int_{B_{R, \sqrt{d} \cdot \bar{h}_0}^+} \sum_{i=1}^d \frac{1}{12} \mathbb{E}_{P_R} \sum_{j=1}^d R_{ij}^2 h_j^2 \left(\frac{\partial f}{\partial x_j} \right)^2 dP_X \\
 &\geq \int_{B_{R, \sqrt{d} \cdot \bar{h}_0}^+ \cap \mathcal{A}_f} \frac{1}{12} h_0^2 c_f^2 dP_X \\
 &\geq \frac{1}{12} \left(\frac{R}{2} \right)^d c_0^2 P_X(\mathcal{A}_f) c_f^2 \cdot \bar{h}_0^2. \tag{B.31}
 \end{aligned}$$

Thus, the assertion is proved. \blacksquare

B.2.3. LOWER BOUND OF SAMPLE ERROR FOR HTR

Proof [of Proposition 15] For any fixed $j \in \mathcal{I}_H$, we define the random variable Z_j by

$$Z_j := \sum_{i=1}^n \mathbf{1}_{A_j}(X_i).$$

Since the random variables $\{\mathbf{1}_{A_j}(X_i)\}_{i=1}^n$ are i.i.d. Bernoulli distributed with parameter $P(X \in A_j)$, elementary probability theory implies that the random variable Z_j is Binomial distributed with parameters n and $P(X \in A_j)$. Therefore, for any $j \in \mathcal{I}_H$, we have

$$\mathbb{E}(Z_j) = n \cdot P(X \in A_j).$$

Moreover, the HTR regressor $f_{D,H}$ can be defined by

$$f_{D,H}(x) = \begin{cases} \frac{\sum_{i=1}^n Y_i \mathbf{1}_{A_j}(X_i)}{\sum_{i=1}^n \mathbf{1}_{A_j}(X_i)} \cdot \mathbf{1}_{A_j}(x) & \text{if } Z_j > 0, \\ 0 & \text{if } Z_j = 0. \end{cases}$$

By the law of total probability, we get

$$\begin{aligned}
 &\mathbb{E}_{P_X} (f_{D,H}(X) - f_{P,H}(X))^2 \\
 &= \sum_{j \in \mathcal{I}_H} \mathbb{E}_{P_X} ((f_{D,H}(X) - f_{P,H}(X))^2 | X \in A_j) \\
 &\quad \cdot P(X \in A_j) \\
 &= \sum_{j \in \mathcal{I}_H} \mathbb{E}_{P_X} ((f_{D,H}(X) - f_{P,H}(X))^2 | X \in A_j, Z_j > 0) \\
 &\quad \cdot P(Z_j > 0) \cdot P(X \in A_j) \tag{B.32} \\
 &+ \sum_{j \in \mathcal{I}_H} \mathbb{E}_{P_X} ((f_{D,H}(X) - f_{P,H}(X))^2 | X \in A_j, Z_j = 0) \\
 &\quad \cdot P(Z_j = 0) \cdot P(X \in A_j). \tag{B.33}
 \end{aligned}$$

For the term (B.32), we have

$$\begin{aligned}
 &\sum_{j \in \mathcal{I}_H} \mathbb{E}_{P_X} ((f_{D,H}(X) - f_{P,H}(X))^2 | X \in A_j, Z_j > 0) \\
 &\quad \cdot P(Z_j > 0) P(X \in A_j)
 \end{aligned}$$

$$\begin{aligned}
 &= \sum_{j \in \mathcal{I}_H} \left(\frac{\sum_{i=1}^n Y_i \mathbf{1}_{A_j}(X_i)}{\sum_{i=1}^n \mathbf{1}_{A_j}(X_i)} - \mathbb{E}(f_{L,P}^*(X) | X \in A_j) \right)^2 \\
 &\quad \cdot P(Z_j > 0) P(X \in A_j) \\
 &= \sum_{j \in \mathcal{I}_H} \left(\sum_{i=1}^n \mathbf{1}_{A_j}(X_i) (Y_i - \mathbb{E}(f_{L,P}^*(X) | X \in A_j)) \right)^2 \\
 &\quad \cdot \frac{P(X \in A_j)}{(\sum_{i=1}^n \mathbf{1}_{A_j}(X_i))^2} \cdot P(Z_j > 0),
 \end{aligned}$$

which yields that for a fixed $j \in \mathcal{I}_H$, there holds

$$\begin{aligned}
 &\mathbb{E} \left(\sum_{j \in \mathcal{I}_H} \left(\sum_{i=1}^n \mathbf{1}_{A_j}(X_i) (Y_i - \mathbb{E}(f_{L,P}^*(X) | X \in A_j)) \right)^2 \right. \\
 &\quad \left. \cdot \frac{P(X \in A_j)}{(\sum_{i=1}^n \mathbf{1}_{A_j}(X_i))^2} \Big| X_i \in A_j \right) \\
 &= \sum_{j \in \mathcal{I}_H} \sum_{i=1}^n \mathbf{1}_{A_j}^2(X_i) \mathbb{E}((Y - f_{P,H}(X))^2 | X \in A_j) \\
 &\quad \cdot \frac{P(X \in A_j)}{(\sum_{i=1}^n \mathbf{1}_{A_j}(X_i))^2} \\
 &= \sum_{j \in \mathcal{I}_H} \frac{P(X \in A_j)}{\sum_{i=1}^n \mathbf{1}_{A_j}(X_i)} \cdot \mathbb{E}((Y - f_{P,H}(X))^2 | X \in A_j). \tag{B.34}
 \end{aligned}$$

Obviously, for any fixed $j \in \mathcal{I}_H$, there holds

$$\mathbb{E}(f_{P,H}(X) | X \in A_j) = \mathbb{E}(f_{L,P}^*(X) | X \in A_j)$$

and consequently we obtain

$$\begin{aligned}
 &\mathbb{E}((Y - f_{P,H}(X))^2 | X \in A_j) \\
 &= \mathbb{E}((Y - f_{L,P}^*(X))^2 | X \in A_j) \\
 &\quad + \mathbb{E}((f_{L,P}^*(X) - f_{P,H}(X))^2 | X \in A_j) \\
 &= \sigma^2 + \mathbb{E}((f_{L,P}^*(X) - f_{P,H}(X))^2 | X \in A_j).
 \end{aligned}$$

Taking expectation over both sides of (B.34) with respect to P^n , we get

$$\begin{aligned}
 &\mathbb{E}_{D \sim P^n} \mathbb{E}_{P_X} (f_{D,H}(X) - f_{P,H}(X))^2 \\
 &= \mathbb{E}_{D \sim P^n} (\mathbb{E}(\mathbb{E}_{P_X} (f_{D,H}(X) - f_{P,H}(X))^2 | X_i \in A_j)) \\
 &= (\sigma^2 + \mathbb{E}(f_{L,P}^*(X) - f_{P,H}(X))^2) \\
 &\quad \cdot \sum_{j \in \mathcal{I}_H} \left(P(X \in A_j) \mathbb{E}_{D \sim P^n} \left(\left(\sum_{i=1}^n \mathbf{1}_{A_j}(X_i) \right)^{-1} \Big| Z_j > 0 \right) \right) \\
 &\quad \cdot P(Z_j > 0) \\
 &= (\sigma^2 + \mathbb{E}(f_{L,P}^*(X) - f_{P,H}(X))^2) \\
 &\quad \cdot \sum_{j \in \mathcal{I}_H} (P(X \in A_j) \mathbb{E}_{D \sim P^n} (Z_j^{-1} | Z_j > 0)) P(Z_j > 0) \\
 &= n^{-1} (\sigma^2 + \mathbb{E}(f_{L,P}^*(X) - f_{P,H}(X))^2)
 \end{aligned}$$

$$\cdot \sum_{j \in \mathcal{I}_H} (\mathbb{E}(Z_j) \cdot \mathbb{E}(Z_j^{-1} | Z_j > 0)) \mathbb{P}(Z_j > 0).$$

Clearly, x^{-1} is convex for $x > 0$. Therefore, by Jensen's inequality, we get

$$\begin{aligned} & \mathbb{E}(Z_j) \cdot \mathbb{E}(Z_j^{-1} | Z_j > 0) \mathbb{P}(Z_j > 0) \\ & \geq \mathbb{E}(Z_j) \cdot \mathbb{E}(Z_j | Z_j > 0)^{-1} \mathbb{P}(Z_j > 0) \\ & = \mathbb{E}(Z) \cdot \mathbb{E}(Z \mathbf{1}_{\{Z > 0\}})^{-1} \mathbb{P}(Z > 0) \mathbb{P}(Z > 0) \\ & = \mathbb{P}(Z > 0)^2 = (1 - \mathbb{P}(Z = 0))^2 \\ & = (1 - (1 - \mathbb{P}(X \in A_j))^n)^2 \\ & \geq 1 - 2e^{-n\mathbb{P}(X \in A_j)}, \end{aligned}$$

where the last inequality follows from $(1 - x)^n \leq e^{-nx}$, $x \in (0, 1)$.

We now turn to estimate the term (B.33). By the definition of $f_{D,H}$, we have

$$\begin{aligned} & \sum_{j \in \mathcal{I}_H} \mathbb{E}_{\mathbb{P}_X} ((f_{D,H}(X) - f_{P,H}(X))^2 | X \in A_j, Z_j = 0) \\ & \quad \cdot \mathbb{P}(Z_j = 0) \cdot \mathbb{P}(X \in A_j) \\ & = \sum_{j \in \mathcal{I}_H} \mathbb{E}_{\mathbb{P}_X} ((f_{P,H}(X))^2 | X \in A_j) \cdot \mathbb{P}(Z_j = 0) \\ & \quad \cdot \mathbb{P}(X \in A_j) \\ & \geq 0. \end{aligned}$$

Let us denote

$$\mathcal{I}_H^{(1)} := \{j \in \mathcal{I}_H : A_j \cap B_R = A_j\}$$

and

$$\mathcal{I}_H^{(2)} := \mathcal{I}_H \setminus \mathcal{I}_H^{(1)}.$$

Then we obviously have $\mathbb{P}(X \in A_j) = \mu(A_j) \geq \underline{h}_0^d$ for all $j \in \mathcal{I}_H^{(1)}$. Combing the above results, we obtain

$$\begin{aligned} & \mathbb{E}_{D \sim \mathbb{P}^n} \mathbb{E}_{\mathbb{P}_X} (f_{D,H}(X) - f_{P,H}(X))^2 \\ & = \sum_{j \in \mathcal{I}_H} \mathbb{E}_{\mathbb{P}_X} ((f_{D,H}(X) - f_{P,H}(X))^2 | X \in A_j, Z_j > 0) \\ & \quad \cdot \mathbb{P}(Z_j > 0) \cdot \mathbb{P}(X \in A_j) \\ & + \sum_{j \in \mathcal{I}_H} \mathbb{E}_{\mathbb{P}_X} ((f_{D,H}(X) - f_{P,H}(X))^2 | X \in A_j, Z_j = 0) \\ & \quad \cdot \mathbb{P}(Z_j = 0) \cdot \mathbb{P}(X \in A_j) \\ & \geq \sum_{j \in \mathcal{I}_H} \mathbb{E}_{\mathbb{P}_X} ((f_{D,H}(X) - f_{P,H}(X))^2 | X \in A_j, Z_j > 0) \\ & \quad \cdot \mathbb{P}(Z_j > 0) \cdot \mathbb{P}(X \in A_j) \\ & = \sum_{j \in \mathcal{I}_H^{(1)}} \mathbb{E}_{\mathbb{P}_X} ((f_{D,H}(X) - f_{P,H}(X))^2 | X \in A_j, Z_j > 0) \end{aligned}$$

$$\begin{aligned} & \cdot \mathbb{P}(Z_j > 0) \cdot \mathbb{P}(X \in A_j) \\ & + \sum_{j \in \mathcal{I}_H^{(2)}} \mathbb{E}_{\mathbb{P}_X} ((f_{D,H}(X) - f_{P,H}(X))^2 | X \in A_j, Z_j > 0) \\ & \quad \cdot \mathbb{P}(Z_j > 0) \cdot \mathbb{P}(X \in A_j) \\ & \geq \sum_{j \in \mathcal{I}_H^{(1)}} \mathbb{E}_{\mathbb{P}_X} ((f_{D,H}(X) - f_{P,H}(X))^2 | X \in A_j, Z_j > 0) \\ & \quad \cdot \mathbb{P}(Z_j > 0) \cdot \mathbb{P}(X \in A_j) \\ & \geq \frac{1}{n} \sum_{j \in \mathcal{I}_H^{(1)}} (1 - 2e^{-n\mathbb{P}(X \in A_j)}) \\ & \quad \cdot (\mathbb{E}(f_{L,P}^*(X) - f_{P,H}(X))^2 + \sigma^2) \\ & \geq \frac{\sigma^2}{n} \left(|\mathcal{I}_H^{(1)}| - \sum_{j \in \mathcal{I}_H^{(1)}} 2e^{-n\mathbb{P}(X \in A_j)} \right). \end{aligned}$$

Therefore, we have

$$\begin{aligned} & \mathbb{E}_{D \sim \mathbb{P}^n} \mathbb{E}_{\mathbb{P}_X} (f_{D,H}(X) - f_{P,H}(X))^2 \\ & \geq \frac{\sigma^2}{n} \left(|\mathcal{I}_H^{(1)}| - \sum_{j \in \mathcal{I}_H^{(1)}} 2e^{-n\mathbb{P}(X \in A_j)} \right) \\ & = \frac{\sigma^2}{n} \left(|\mathcal{I}_H^{(1)}| - 2|\mathcal{I}_H^{(1)}| \exp(-n\underline{h}_0^d) \right) \\ & \geq \frac{\sigma^2}{n} \left(\frac{2R - \sqrt{d} \cdot \bar{h}_0}{\bar{h}_0} \right)^d \left(1 - \frac{2}{e} \right) \\ & \geq 4R^d \sigma^2 (1 - 2e^{-1}) \bar{h}_0^{-d} n^{-1}, \end{aligned} \quad (\text{B.35})$$

where the last inequality follows from Assumption 1. \blacksquare

B.2.4. PROOF RELATED TO SECTION 3.2

Proof [of Theorem 2] Theorem 11 together with Proposition 13 implies

$$\begin{aligned} & \mathcal{R}_{L_{\bar{h}_0}, \mathbb{P}}(f_{D,B}) - \mathcal{R}_{L_{\bar{h}_0}, \mathbb{P}}^* \\ & \lesssim \lambda_1/T + \lambda_2 \underline{h}_0^{-2d} + \bar{h}_0^{2(1+\alpha)} + T^{-1} \bar{h}_0^2 \\ & \quad + (T/\lambda_1)^{2\delta'} \lambda_2^{-1} n^{-2}, \end{aligned}$$

where $\delta' := 1 - \delta$ and $\delta := (\underline{h}_0/c_d)^d$. Choosing

$$\begin{aligned} \lambda_{1,n} & := n^{-\frac{2}{2(1+\alpha)(2-\delta)+d}}, \\ \lambda_{2,n} & := n^{-\frac{2(\alpha+d+1)}{2(1+\alpha)(2-\delta)+d}}, \\ \bar{h}_{0,n} & := n^{-\frac{1}{2(1+\alpha)(2-\delta)+d}}, \\ T_n & := n^{\frac{2\alpha}{2(1+\alpha)(2-\delta)+d}}, \end{aligned}$$

we obtain

$$\mathcal{R}_{L_{\bar{h}_0}, \mathbb{P}}(f_{D,B}) - \mathcal{R}_{L_{\bar{h}_0}, \mathbb{P}}^* \lesssim n^{-\frac{2(1+\alpha)}{2(1+\alpha)(2-\delta)+d}}.$$

with probability P^n not less than $1 - 3e^{-\tau}$ in expectation with respect to P_H . This completes the proof. ■

Proof [of Theorem 3] Recall the error decomposition (A.11). Using the estimates (B.31) and (B.35) and choosing $\bar{h}_{0,n} := n^{-\frac{1}{d+2}}$, we get

$$\begin{aligned} & \mathbb{E}_{\nu_n}(\mathcal{R}_{L,P}(f_{D,H_n}) - \mathcal{R}_{L,P}^*) \\ &= \mathbb{E}_{\nu_n} \mathbb{E}_{P_X}(f_{D,H_n}(X) - f_{L,P}^*(X))^2 \\ &\geq \frac{d}{12} \left(\frac{R}{2}\right)^d c_0^2 P_X(\mathcal{A}_f) \underline{c}_f^2 \cdot \bar{h}_{0,n}^2 + \frac{4R^2\sigma^2}{n} (1 - 2e^{-1}) \bar{h}_{0,n}^{-d} \\ &\gtrsim n^{-\frac{2}{2+d}}. \end{aligned}$$

Consequently we have

$$\mathbb{E}_{\nu_n}(\mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P}^*) \gtrsim n^{-\frac{2}{2+d}},$$

which proves the assertion. ■

C. Description of Real Datasets

ABA: The *Abalone* dataset originally comes from a biology study (Tarbath, 2003) which focused on the relationship between the age of an abalone and its other features. Now accessible on UCI datasets, abalone dataset consists 4,177 observations on abalone ages with 8 attributes, including physical measurements on abalones and the environment.

BOD: The *Body-fat* dataset is available on *Libstat* of CMU, containing 252 observations with 13 attributes of physical measurement of human body, and each observation contains two targeted variables, body density and the percentage of body fat.

HOU: The *Housing-Boston* dataset can be acquired from *LibSVM* datasets of NTU, which is comprised of 506 observations with 13 features. The dataset is used to predict the price of an house in Boston.

MG: This dataset can be traced back to Flake & Lawrence (2002). It consists of 1,385 observations of dimension 6.

MPG: The *Auto MPG* dataset is a modified version of MPG dataset in *Libstat* of CMU, containing 398 instances with 8 attributes. Compared with its original version, 8 instances are deleted since their missing *mpg* target value.

PYR: The *Pyrimidines* dataset is a subset of *Qualitative Structure Activity Relationships* dataset on UCI datasets. This sub-dataset has 74 instances of dimension 27.

SPA: The *Geographical Analysis Spatial* dataset is accessible in *Libstat* of CMU, originally uploaded by Pace & Barry (1997). It comprises 3,107 observations of dimension 6.

TRI: The *Triazines* dataset is another subset of *Qualitative Structure Activity Relationships* dataset on UCI datasets. This part consists 186 instances of dimension 60.

References

- Breiman, L. Some infinity theory for predictor ensembles. Technical report, Technical Report 579, Statistics Dept. UCB, 2000.
- Flake, G. W. and Lawrence, S. Efficient SVM regression training with SMO. *Machine Learning*, 46(1-3):271–290, 2002.
- Pace, R. K. and Barry, R. Quick computation of regressions with a spatially autoregressive dependent variable. *Geographical Analysis*, 29(3):232–247, 1997.
- Steinwart, I. and Christmann, A. *Support Vector Machines*. Information Science and Statistics. Springer, New York, 2008.
- Tarbath, D. Population parameters of blacklip abalone (*halotis rubra* leach) at the actaeons in south-east tasmania. 2003. URL <https://academic.microsoft.com/paper/168232552>.
- Van der Vaart, A. W. and Wellner, J. A. *Weak Convergence and Empirical Processes*. Springer Series in Statistics. Springer-Verlag, New York, 1996.