



HAL
open science

Minimax adaptive dimension reduction for regression

Quentin Paris

► **To cite this version:**

| Quentin Paris. Minimax adaptive dimension reduction for regression. 2012. hal-00768911v1

HAL Id: hal-00768911

<https://hal.science/hal-00768911v1>

Preprint submitted on 27 Dec 2012 (v1), last revised 28 Jun 2013 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Minimax adaptive dimension reduction for regression

Quentin PARIS

IRMAR, ENS Cachan Bretagne, CNRS, UEB
Campus de Ker Lann
Avenue Robert Schuman, 35170 Bruz, France
quentin.paris@bretagne.ens-cachan.fr

Abstract

Let (X, Y) be a random variable taking its values in $\mathcal{X} \times \mathbb{R}$, where $\mathcal{X} \subset \mathbb{R}^p$. We study the estimation of the regression function $r(x) := \mathbb{E}(Y|X = x)$ from a minimax point of view, assuming that r belongs to a class of functions \mathcal{F} of the form

$$\mathcal{G} \circ \mathcal{H} := \left\{ g \circ h : g \in \mathcal{G}, h \in \mathcal{H} \right\}.$$

Here, \mathcal{G} is a class of β -Hölder functions $g : \mathbb{R}^p \rightarrow \mathbb{R}$ and \mathcal{H} a class of functions $h : \mathcal{X} \rightarrow \mathbb{R}^p$ which need not be linear nor regular. We define the *reduced dimension* as the smallest ℓ for which $r = g \circ h$ for some $g \in \mathcal{G}$ and some $h \in \mathcal{H}$ such that $h(\mathcal{X})$ spans a subspace of dimension ℓ . In this context, we construct an adaptive estimate which converges at a rate depending only on the *reduced dimension* under entropy conditions on \mathcal{H} .

Index Terms – Regression estimation, dimension reduction, minimax rates of convergence, empirical risk minimization, metric entropy.

AMS 2000 Classification – 62H12, 62G08.

1 Introduction

From a general point of view, the goal of regression is to infer about the conditional distribution of a real-valued response variable Y given an \mathcal{X} -valued predictor variable X where $\mathcal{X} \subset \mathbb{R}^p$. In the statistical framework, one usually focuses on the estimation of the regression function

$$r(x) := \mathbb{E}(Y|X = x), \quad (1.1)$$

based on a sample $(X_1, Y_1), \dots, (X_n, Y_n)$ of n independent and identically distributed random variables with same distribution P as (X, Y) .

When addressing this task, the statistician considers a set of assumptions on the underlying distribution of the observations based on some physical understanding of the phenomenon under study. Roughly speaking, this modelization step consists in specifying a class \mathcal{D} of distributions of the random variable (X, Y) for which the regression function r belongs to some class \mathcal{F} of functions $f : \mathcal{X} \rightarrow \mathbb{R}$. Given such a regression model, one usually refers to optimal rates of convergence in the minimax sense as benchmarks to assess the performance of a particular estimate. In this context, optimal rates of convergence may be defined as follows. The sequence (v_n) is said to be an optimal rate of convergence in the minimax sense if it is a lower minimax rate, *i.e.* if

$$\liminf_{n \rightarrow +\infty} v_n^{-2} \inf_{\hat{r}} \sup_{P \in \mathcal{D}} \mathbb{E}_{(X,Y) \sim P} (\hat{r}(X) - \mathbb{E}(Y|X))^2 > 0, \quad (1.2)$$

where the infimum is taken over all estimates \hat{r} based on our sample, and if there exists an estimate \hat{r} such that

$$\limsup_{n \rightarrow +\infty} v_n^{-2} \sup_{P \in \mathcal{D}} \mathbb{E}_{(X,Y) \sim P} (\hat{r}(X) - \mathbb{E}(Y|X))^2 < +\infty. \quad (1.3)$$

In (1.2) and (1.3), we have used the subscript $(X, Y) \sim P$ to stress that (X, Y) is drawn from P . When no confusion may arise, we will omit it in the rest of the paper.

A major problem, known as the *curse of dimensionality*, is that for most models used across the literature, the optimal rates of convergence depend on the dimension p of the predictor X . For instance, suppose \mathcal{F} is the class of all $[0, 1]$ -valued, β -Hölder functions defined on the open unit Euclidian ball in \mathbb{R}^p . Then

it can be shown that, under some technical conditions, the optimal rate of convergence in this context is $n^{-\beta/(2\beta+p)}$ (see *e.g.* Theorem 3.2 in Györfi *et al.* (2002) and Theorem 1 in Kolher *et al.* (2009)). Hence, when the dimension p is large, the optimal rate of convergence is slow. To face this drawback when dealing with high-dimensional data, one usually considers a model which encodes so called *structural assumptions* in addition to reasonable regularity conditions, thus reducing the overall complexity of the model. Many such models have been studied and proved effective in practice among which we mention additive regression models, projection pursuit, single index models (see *e.g.* Chapter 22 in the book by Györfi *et al.* (2002) and the references cited therein) and the body of theory and methods known as sufficient dimension reduction (see *e.g.* the foundational article of Li (1991) which introduces the SIR method or the article of Cook and Li (2002) which defines the important notion of central mean subspace).

In the present article we study a regression model with a structural assumption inspired by sufficient dimension reduction methodology. Precisely, we assume that the regression function belongs to a class \mathcal{F} of the form

$$\mathcal{G} \circ \mathcal{H} := \{g \circ h : g \in \mathcal{G}, h \in \mathcal{H}\}, \quad (1.4)$$

where \mathcal{G} denotes a class of functions $g : \mathbb{R}^p \rightarrow \mathbb{R}$ and where \mathcal{H} denotes a class of functions $h : \mathcal{X} \rightarrow \mathbb{R}^p$. In this context, having $r \in \mathcal{F}$ means that the two conditions

$$(i) \mathbb{E}(Y|h(X)) = \mathbb{E}(Y|X), \quad \text{and} \quad (ii) \mathbb{E}(Y|h(X) = \cdot) \in \mathcal{G},$$

are satisfied for at least one function $h \in \mathcal{H}$. From a heuristic point of view, the main interest here is that condition (i) should allow to replace the high dimensional predictor X by a possibly low dimensional predictor $h(X)$. In this approach, it is expected that the rates of convergence depend on the dimension of $h(X)$ rather than on the dimension of X . Note that condition (i) generalizes the usual assumption made in sufficient dimension reduction where one assumes that h is a matrix. In the present paper the functions in \mathcal{H} need not be linear nor regular. The class \mathcal{G} is a class of regular functions the construction of which is detailed in the next section. In the classical sufficient dimension reduction context, a similar approach has been successfully applied by Cadre and Dong (2010) to obtain fast rates of convergence.

From a technical point of view, the impact of dimensionality on rates of convergence can be explained as follows. A first important fact is that optimal rates

of convergence usually depend on a measure of the complexity of \mathcal{F} . Suppose that the marginal distribution μ of X is fixed and that the complexity of \mathcal{F} is measured in terms of the ε -metric entropy of \mathcal{F} in $\mathbb{L}^2(\mu)$ denoted $H(\varepsilon, \mathcal{F}, \mu)$ (i.e. the logarithm of the minimal number of $\mathbb{L}^2(\mu)$ -balls of radius ε that are needed to cover \mathcal{F}). Then, if there exists $c > 0$ such that $H(\varepsilon, \mathcal{F}, \mu) \geq c\varepsilon^{-s}$, the optimal rates of convergence are lower bounded by $n^{-1/(2+s)}$ under some technical assumptions (see *e.g.* Theorem 6 in Yang and Barron (1999)). Suppose this time the complexity of \mathcal{F} is measured in terms of $H(\varepsilon, \mathcal{F}) := \sup H(\varepsilon, \mathcal{F}, Q)$, where the supremum is taken over all probability measures with finite support in \mathcal{X} . Then if there exists $C > 0$ such that $H(\varepsilon, \mathcal{F}) \leq C\varepsilon^{-s}$, the optimal rates of convergence are upper bounded by $n^{-1/(2+s)}$ under some technical assumptions (see Appendix A). The second important fact is that the value of the exponent s usually depends on the dimension p of the predictor X . For example, when \mathcal{F} is the class of all $[0, 1]$ -valued, β -Hölder functions defined on the open unit Euclidian ball in \mathbb{R}^p , the typical value for the exponent s is p/β which explains that the optimal rate in this context is $n^{-1/(2+s)} = n^{-\beta/(2\beta+p)}$. In the case of a model as defined in (1.4), we will use the composite structure to get bounds on the entropy of \mathcal{F} which depend only the regularity of the functions in \mathcal{G} and the entropy of \mathcal{H} .

The paper is organized as follows. In Section 2, we describe precisely our model and the assumptions made on the classes \mathcal{G} and \mathcal{H} . In Section 3, we fix an integer $\ell \in \{1, \dots, p\}$ and derive bounds for the optimal rates in the context of our composite model when all functions $h \in \mathcal{H}$ satisfy $\dim S(h) \leq \ell$ where $S(h)$ denotes the subspace of \mathbb{R}^p spanned by $h(\mathcal{X})$. In Section 4 we present our main result, namely Theorem 4.3, which studies the general case where no restriction is made on the dimension of $S(h)$ for $h \in \mathcal{H}$. In this case, using the results of Section 3, we construct an estimate \hat{r} of the regression function which adapts to the *reduced dimension* d defined as the smallest ℓ for which one has $r = g \circ h$ for some $g \in \mathcal{G}$ and some $h \in \mathcal{H}$ with $\dim S(h) \leq \ell$. More precisely, we prove that uniformly over a certain class of distributions, our estimate converges to r at a rate depending only on the reduced dimension d . Furthermore, this rate corresponds to the rate one would obtain if the true value of the reduced dimension were known as exposed in Section 3. We have reported in the appendix a useful result concerning rates of convergence for least-squares estimates in the case of bounded regression. The proof follows from the lines devised in Koltchinskii (2006) and is available in the supplementary material.

2 Definition of the model

Let (X, Y) be an $\mathcal{X} \times \mathbb{R}$ -valued random variable of distribution P where $\mathcal{X} \subset \mathbb{R}^p$. We denote μ the distribution of X and for $x \in \mathcal{X}$ we set

$$r(x) := \mathbb{E}(Y|X = x). \quad (2.1)$$

Let \mathcal{F} be a class of functions $f : \mathcal{X} \rightarrow \mathbb{R}$ of the form

$$\mathcal{F} := \mathcal{G} \circ \mathcal{H} := \left\{ g \circ h : g \in \mathcal{G}, h \in \mathcal{H} \right\},$$

where \mathcal{G} and \mathcal{H} are taken as follows. Let $R > 0$. For all $h \in \mathcal{H}$ we assume that

$$\|h\|_{\mathcal{X}} := \sup_{x \in \mathcal{X}} \|h(x)\| < R, \quad (2.2)$$

where $\|\cdot\|$ stands for the Euclidean norm in \mathbb{R}^p . Let \mathbf{B} be the open Euclidean ball in \mathbb{R}^p with center the origin and radius R . Let $\beta > 0$. We denote $\mathcal{C}^{[\beta]}(\mathbf{B})$ the space of $[\beta]$ -times continuously differentiable functions from \mathbf{B} to \mathbb{R} , where $[\beta]$ stands for the greatest integer strictly smaller than β . For all g in $\mathcal{C}^{[\beta]}(\mathbf{B})$ we set

$$\|g\|_{\beta} := \max_{|s| \leq [\beta]} \|\partial^s g\|_{\infty} + \max_{|s| = [\beta]} \sup_{u \neq u'} \frac{|\partial^s g(u) - \partial^s g(u')|}{\|u - u'\|^{\beta - [\beta]}}, \quad (2.3)$$

where for every multi-index $s = (s_1, \dots, s_p) \in \mathbb{N}^p$, we have denoted $|s| := \sum_i s_i$ and $\partial^s := \partial_1^{s_1} \dots \partial_p^{s_p}$. For $L > 0$, the class \mathcal{G} is defined by

$$\mathcal{G} := \left\{ g \in \mathcal{C}^{[\beta]}(\mathbf{B}) : \|g\|_{\beta} \leq L \right\}. \quad (2.4)$$

In the sequel, we suppose that the distribution P of (X, Y) belongs to the following class of distributions.

Class \mathcal{D} of distributions. Let $\tau > 0$ and $B > 0$ be fixed. Let the marginal distribution μ be fixed. We denote \mathcal{D} the class of distributions P of the random variable (X, Y) such that X is of distribution μ , such that the regression function r belongs to \mathcal{F} and such that Y satisfies the exponential moment condition

$$\mathbb{E} e^{\tau|Y|} \leq B. \quad (2.5)$$

3 Non adaptive case

Fix $\ell \in \{1, \dots, p\}$. In this section, we derive bounds on the optimal rates of convergence associated to the class of distributions \mathcal{D} under the assumption that all functions $h \in \mathcal{H}$ satisfy

$$\dim S(h) \leq \ell, \quad (3.1)$$

where $S(h)$ denotes the subspace of \mathbb{R}^p spanned by $h(\mathcal{X})$. Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be a sample of n independent and identically distributed random variable with same distribution P as (X, Y) .

First we present an upper bound. For any Borel measure Q on \mathcal{X} , we denote by $\mathbb{L}^2(Q)$ the set of measurable functions $f : \mathcal{X} \rightarrow \mathbb{R}$ such that

$$\|f\|_Q^2 := \int_{\mathcal{X}} |f|^2 dQ < +\infty. \quad (3.2)$$

For a class $\mathcal{C} \subset \mathbb{L}^2(Q)$, we recall that the ε -covering number of \mathcal{C} with respect to $\|\cdot\|_Q$, denoted $N(\varepsilon, \mathcal{C}, Q)$, is the minimal number of $\|\cdot\|_Q$ -balls of radius ε that are needed to cover \mathcal{C} . The ε -metric entropy of \mathcal{C} with respect to $\|\cdot\|_Q$ is defined by $H(\varepsilon, \mathcal{C}, Q) := \ln N(\varepsilon, \mathcal{C}, Q)$. In this paper, we set

$$H(\varepsilon, \mathcal{C}) := \sup_Q H(\varepsilon, \mathcal{C}, Q), \quad (3.3)$$

where the supremum is taken over all probability measures Q with finite support in \mathcal{X} . These notations will also be used for classes of vector valued functions. In this case, it will be understood that the absolute value $|\cdot|$ in (3.2) is replaced by the appropriate Euclidean norm.

We fix a constant $\alpha > 2$ and define the least-squares estimate \hat{f} as any element in \mathcal{F} satisfying

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \left(Y_i \mathbf{1}\{|Y_i| \leq T_n\} - f(X_i) \right)^2, \quad (3.4)$$

where $T_n := (\ln n)^{\alpha/2}$. Implicitly, it will be assumed that such an element exists, though it might not be unique. We denote

$$v_{n,\ell} := \left(\frac{(\ln n)^\alpha}{n} \right)^{\beta/(2\beta+\ell)}.$$

Theorem 3.1. *Suppose that $\beta \geq 1$, that $\beta > \ell/2$ and suppose that there exists $C > 0$ and $0 < s \leq \ell/\beta$ such that for all $\varepsilon > 0$ we have*

$$H(\varepsilon, \mathcal{H}) \leq C\varepsilon^{-s}. \quad (3.5)$$

Then we have

$$\limsup_{n \rightarrow +\infty} v_{n,\ell}^{-2} \sup_{P \in \mathcal{D}} \mathbb{E} (\hat{f}(X) - \mathbb{E}(Y|X))^2 < +\infty. \quad (3.6)$$

An immediate consequence of Theorem 3.1 is that, up to a logarithmic factor, the optimal rate of convergence associated to \mathcal{D} is upper bounded by $n^{-\beta/(2\beta+\ell)}$ which depends only on ℓ and not anymore on the dimension p of the predictor X . Note that, under some technical conditions, $n^{-\beta/(2\beta+\ell)}$ corresponds to the optimal rate of convergence for the estimation of a β -Hölder function in the case where X is ℓ -dimensional (see *e.g.* Theorem 3.2 in Györfi *et al.* (2002) and Theorem 1 in Kolher *et al.* (2009)). If the exponential moment condition (2.5) is replaced by a boundedness assumption on Y , then a slight modification of the proof reveals that Theorem 3.1 holds with $v_{n,\ell}$ replaced by $n^{-\beta/(2\beta+\ell)}$.

Example 3.2. *An example where condition (3.5) is fulfilled is the following. Consider the case where \mathcal{H} is a parametric class of the form*

$$\mathcal{H} = \left\{ h_\theta : \theta \in \Theta \subset \mathbb{R}^k \right\}, \quad (3.7)$$

where Θ is bounded. Suppose that there exists a constant $C > 0$ such that for all $\theta, \theta' \in \Theta$ we have

$$\|h_\theta - h_{\theta'}\|_{\mathcal{X}} \leq C|\theta - \theta'|, \quad (3.8)$$

where $|\cdot|$ stands for the Euclidian norm in \mathbb{R}^k . Then, for all $\varepsilon > 0$ we have $H(\varepsilon, \mathcal{H}) \leq H(\varepsilon/C, \Theta)$ where $H(\varepsilon, \Theta)$ stands for the logarithm of the minimal number of Euclidean balls of radius ε that are needed to cover Θ . Since Θ is bounded, it is included in a Euclidean ball of radius ρ for some $\rho > 0$. Therefore, it follows from Proposition 5 in Cucker and Smale (2001) that for all $\varepsilon > 0$ we have

$$H(\varepsilon, \Theta) \leq k \ln(4\rho/\varepsilon).$$

As a result, there exists a constant C' such that for all $\varepsilon > 0$ we have

$$H(\varepsilon, \mathcal{H}) \leq C' \ln(1/\varepsilon),$$

and condition (3.5) is satisfied.

Now we give a lower bound for the optimal rates of convergence. The proof of this result uses arguments from Yang and Barron (1999).

Theorem 3.3. *Suppose that $\beta > 0$. Suppose that there exists $h \in \mathcal{H}$ with $\dim S(h) = \ell$ and a constant $c > 0$ such that*

$$\mu \circ h^{-1} \geq c \lambda_h(\cdot \cap \mathbf{B}), \quad (3.9)$$

where λ_h denotes the Lebesgue measure in $S(h)$. Then, we have

$$\liminf_{n \rightarrow +\infty} n^{2\beta/(2\beta+\ell)} \inf_{\hat{f}} \sup_{P \in \mathcal{D}} \mathbb{E} (\hat{f}(X) - \mathbb{E}(Y|X))^2 > 0,$$

where the infimum is taken over all estimates \hat{f} .

The condition in (3.9) is not restrictive. As an example, it is satisfied if $\mathcal{X} = \mathbf{B}$, if the function $h : (x_1, \dots, x_p) \in \mathbb{R}^p \mapsto (x_1, \dots, x_\ell, 0, \dots, 0) \in \mathbb{R}^p$ belongs to \mathcal{H} and if μ has a density with respect to the Lebesgue measure which is lower bounded by a positive constant on \mathbf{B} . Theorem 3.3 shows that the optimal rate associated to \mathcal{D} is lower bounded by $n^{-\beta/(2\beta+\ell)}$ which, up to a logarithmic factor, corresponds to the upper bound on the optimal rate given in Theorem 3.1.

4 Dimension adaptivity

In Section 3 we have studied the case where all functions $h \in \mathcal{H}$ satisfy $\dim S(h) \leq \ell$ for some $\ell \in \{1, \dots, p\}$. In this section, we impose no restriction on the dimension of $S(h)$ for $h \in \mathcal{H}$. In Theorem 4.1 we study an estimate of the *reduced dimension* defined hereafter. Based on this estimate, we construct an estimate of the regression function and show in our main result, Theorem 4.3, that it adapts to the reduced dimension.

First we need some notations. For all $\ell \in \{1, \dots, p\}$, let \mathcal{H}_ℓ be the class of all functions $h \in \mathcal{H}$ such that $S(h)$ is at most ℓ -dimensional, *i.e.*

$$\mathcal{H}_\ell := \left\{ h \in \mathcal{H} : \dim S(h) \leq \ell \right\}. \quad (4.1)$$

We define the subset \mathcal{F}_ℓ of \mathcal{F} by

$$\mathcal{F}_\ell := \left\{ g \circ h : g \in \mathcal{G}, h \in \mathcal{H}_\ell \right\}. \quad (4.2)$$

The \mathcal{F}_ℓ 's form a nested family of models, that is $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots \subset \mathcal{F}_p = \mathcal{F}$. The *reduced dimension* associated to the underlying distribution $P \in \mathcal{D}$ of our observations is defined by

$$d_P := \min \left\{ \ell : r \in \mathcal{F}_\ell \right\}. \quad (4.3)$$

When no confusion may arise, we will drop the dependency on P in the notation. Our first task will be to derive a tractable representation of the reduced dimension, suitable for estimation purposes. We shall use the following assumption.

Assumption (A). For all $\ell \in \{1, \dots, p\}$, the set \mathcal{F}_ℓ is compact in $\mathbb{L}^2(\mu)$.

Let R_ℓ be the risk defined by

$$R_\ell := \inf_{f \in \mathcal{F}_\ell} \mathbb{E}(Y - f(X))^2. \quad (4.4)$$

Since the \mathcal{F}_ℓ 's are nested, the function $\ell \in \{1, \dots, p\} \mapsto R_\ell$ is non increasing. Then, using Assumption (A), we deduce that

$$d = \min \left\{ \ell : R_\ell = R_p \right\}. \quad (4.5)$$

Consequently, for all $0 < \delta < \Delta$, we have

$$d = \min \left\{ \ell = 1, \dots, p : R_\ell \leq R_p + \delta \right\}, \quad (4.6)$$

where Δ is defined as

$$\Delta := \min \left\{ R_\ell - R_p : R_\ell > R_p \right\}, \quad (4.7)$$

with the convention $\min \emptyset = +\infty$. Observe that $\Delta > 0$ and that, when $d \geq 2$, Δ corresponds to the distance from r to \mathcal{F}_{d-1} in $\mathbb{L}^2(\mu)$, that is

$$\Delta = \inf_{f \in \mathcal{F}_{d-1}} \|f - r\|_\mu^2. \quad (4.8)$$

Equations (4.5) and (4.8) are proved in Appendix B.

Based on (4.6), one may construct a natural estimate of d . First we introduce the empirical counterpart of the risk R_ℓ defined in (4.4) by setting

$$\hat{R}_\ell := \inf_{f \in \mathcal{F}_\ell} \frac{1}{n} \sum_{i=1}^n \left(Y_i \mathbf{1}\{|Y_i| \leq a_n\} - f(X_i) \right)^2,$$

where (a_n) denotes an increasing sequence of positive numbers. Then, given a decreasing sequence of positive numbers (b_n) , we define the estimate

$$\hat{d} := \min \left\{ \ell = 1, \dots, p : \hat{R}_\ell \leq \hat{R}_p + b_n \right\}. \quad (4.9)$$

The asymptotic behavior of \hat{d} is given by the following result.

Theorem 4.1. *Suppose that:*

- (1) $\beta \geq 1$,
- (2) *There exists $C > 0$ and $0 < s \leq p/\beta$ such that $H(\varepsilon, \mathcal{H}) \leq C\varepsilon^{-s}$ for all $\varepsilon > 0$,*
- (3) *Assumption (A) is satisfied.*

Suppose that $a_n = n^u$ and $b_n = n^{-u'}$ where $u > 0$, where $u' \geq 0$ and where

$$0 < 2u + u' < \frac{2\beta}{2\beta + p}.$$

Then the following two statements hold:

(i) *For all $\vartheta > 0$ we have*

$$\lim_{n \rightarrow +\infty} n^\vartheta \sup_{P \in \mathcal{D}} \mathbb{P}(\hat{d} > d) = 0.$$

(ii) *For all $\underline{\delta} > 0$ and for all $\vartheta > 0$ we have*

$$\lim_{n \rightarrow +\infty} n^\vartheta \sup_{P \in \mathcal{D}(\underline{\delta})} \mathbb{P}(\hat{d} < d) = 0,$$

where we have denoted $\mathcal{D}(\underline{\delta}) := \{P \in \mathcal{D} : \Delta \geq \underline{\delta}\}$.

Theorem 4.1 reveals that for all $\underline{\delta} > 0$, the estimate \hat{d} converges exponentially fast to d , uniformly for all $P \in \mathcal{D}(\underline{\delta})$. This is not the case for $\underline{\delta} = 0$, as shown by the following theorem.

Theorem 4.2. *Let $\mathcal{D}_0 \subset \mathcal{D}$ be an arbitrary subset of \mathcal{D} such that*

$$\inf \{ \Delta : P \in \mathcal{D}_0 \} = 0.$$

Under the conditions of Theorem 4.1 and for all $\vartheta > 0$ we have

$$\sup_{P \in \mathcal{D}_0} \mathbb{P}(\hat{d} < d) \geq 1 - \frac{1}{n^\vartheta},$$

provided n is large enough.

Now we apply these results to the estimation of the regression function. First, for all $\ell \in \{1, \dots, p\}$, let \hat{f}_ℓ be defined as any element in \mathcal{F}_ℓ satisfying

$$\hat{f}_\ell \in \arg \min_{f \in \mathcal{F}_\ell} \frac{1}{n} \sum_{i=1}^n \left(Y_i \mathbf{1}\{|Y_i| \leq T_n\} - f(X_i) \right)^2, \quad (4.10)$$

where T_n is defined as in (3.4). Our next result proves that the estimate \hat{f} defined by

$$\hat{f} := \hat{f}_{\hat{d}}, \quad (4.11)$$

adapts to the reduced dimension d in the sense that for all $\underline{\delta} > 0$ and uniformly for all $P \in \mathcal{D}(\underline{\delta})$ it converges to r at a rate depending only on d .

Theorem 4.3. *Suppose that:*

- (1') $\beta \geq 1$ and $\beta > p/2$,
- (2') For all $\ell \in \{1, \dots, p\}$, there exists $C > 0$ and $0 < s \leq \ell/\beta$ such that for all $\varepsilon > 0$ we have $H(\varepsilon, \mathcal{H}_\ell) \leq C\varepsilon^{-s}$,
- (3) Assumption **(A)** is satisfied.

Suppose that $a_n = n^u$ and $b_n = n^{-u'}$ where $u > 0$, where $u' \geq 0$ and where

$$0 < 2u + u' < \frac{2\beta}{2\beta + p}.$$

Then for all $\underline{\delta} > 0$ we have

$$\limsup_{n \rightarrow +\infty} \sup_{P \in \mathcal{D}(\underline{\delta})} \mathbf{v}_{n,d}^{-2} \mathbb{E} \left(\hat{f}(X) - \mathbb{E}(Y|X) \right)^2 < +\infty,$$

where it is understood that $d = d_p$.

Theorem 4.3 shows that the estimate \hat{f} defined in (4.11) converges to r as fast as the estimate \hat{f}_d one would choose if the true value of the reduced dimension d were known as exposed in Section 3. Therefore, for all $\underline{\delta} > 0$, the estimate \hat{f} adapts to the reduced dimension d , uniformly over all $P \in \mathcal{D}(\underline{\delta})$. Note that for condition (2') to be satisfied for all $\ell \in \{1, \dots, p\}$, it is sufficient to impose that $H(\varepsilon, \mathcal{H}) \leq C\varepsilon^{-s}$ for some $C > 0$ and some $0 < s \leq 1/\beta$. For example, this condition is satisfied when \mathcal{H} is chosen as in Example (3.2).

The condition $\beta > p/2$ in Theorem 4.3 allows adaptation for all values of $d \in \{1, \dots, p\}$. We do not know whether this result holds for $\beta \leq p/2$. That being

said, if adaptation is required only for small dimensions we have the following result. For all $\beta \geq 1$ and all $\underline{\delta} > 0$, by considering the smaller class

$$\mathcal{D}(\underline{\delta}, \beta) := \mathcal{D}(\underline{\delta}) \cap \{P \in \mathcal{D} : d \leq \ell_\beta\},$$

where

$$\ell_\beta := [2\beta] - 1,$$

and where $[x]$ stands for the greatest integer smaller or equal to x , we obtain readily that

$$\limsup_{n \rightarrow +\infty} \sup_{P \in \mathcal{D}(\underline{\delta}, \beta)} v_{n,d}^{-2} \mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 < +\infty.$$

Finally, note that if the exponential moment condition (2.5) is replaced by a boundedness assumption on Y , then Theorem 4.3 holds with $v_{n,d}$ replaced by $n^{-\beta/(2\beta+d)}$.

5 Proofs

5.1 Proof of Theorem 3.1

Lemma 5.1. *Assume $\beta \geq 1$. Then for all $\varepsilon > 0$ we have*

$$H(\varepsilon, \mathcal{F}) \leq \sup_{h \in \mathcal{H}} H\left(\frac{\varepsilon}{2}, \mathcal{G} \circ h\right) + H\left(\frac{\varepsilon}{2L}, \mathcal{H}\right).$$

Proof – We fix $\varepsilon > 0$ and let Q be any probability measure with support in \mathcal{X} . We denote

$$N := N\left(\frac{\varepsilon}{2L}, \mathcal{H}, Q\right),$$

and choose an $\frac{\varepsilon}{2L}$ -covering of \mathcal{H} with respect to $\|\cdot\|_Q$ of minimum cardinality $\{h_1, \dots, h_N\}$. For all $i \in \{1, \dots, N\}$, let

$$N_i := N\left(\frac{\varepsilon}{2}, \mathcal{G} \circ h_i, Q\right),$$

and let $\{g_1^i \circ h_i, \dots, g_{N_i}^i \circ h_i\}$ be an $\frac{\varepsilon}{2}$ -covering of $\mathcal{G} \circ h_i$ with respect to $\|\cdot\|_Q$. Then, let $g \in \mathcal{G}$ and $h \in \mathcal{H}$ be chosen arbitrarily. By definition, there exists $i \in \{1, \dots, N\}$ and $j \in \{1, \dots, N_i\}$ such that

$$\sqrt{\int \|h - h_i\|^2 dQ} \leq \frac{\varepsilon}{2L} \quad \text{and} \quad \sqrt{\int (g \circ h_i - g_j^i \circ h_i)^2 dQ} \leq \frac{\varepsilon}{2}.$$

Therefore, we have

$$\begin{aligned}
& \sqrt{\int (g \circ h - g_j^i \circ h_i)^2 dQ} \\
& \leq \sqrt{\int (g \circ h - g \circ h_i)^2 dQ} + \sqrt{\int (g \circ h_i - g_j^i \circ h_i)^2 dQ} \\
& \leq \sqrt{\int (g \circ h - g \circ h_i)^2 dQ} + \frac{\varepsilon}{2}.
\end{aligned}$$

According to the mean value Theorem, all $g \in \mathcal{G}$ are L -Lipschitz. Therefore, we have

$$\begin{aligned}
\sqrt{\int (g \circ h - g_j^i \circ h_i)^2 dQ} & \leq \sqrt{\int (g \circ h - g \circ h_i)^2 dQ} + \frac{\varepsilon}{2} \\
& \leq L \sqrt{\int \|h - h_i\|^2 dQ} + \frac{\varepsilon}{2} \\
& \leq \varepsilon.
\end{aligned} \tag{5.1}$$

As a result, we have obtained that

$$\begin{aligned}
N(\varepsilon, \mathcal{F}, Q) & \leq \left| \{g_j^i \circ h_i : i \in \{1, \dots, N\}, j \in \{1, \dots, N_i\}\} \right| \\
& = \sum_{i=1}^N N_i \\
& \leq \sup_{h \in \mathcal{H}} N\left(\frac{\varepsilon}{2}, \mathcal{G} \circ h, Q\right) N\left(\frac{\varepsilon}{2L}, \mathcal{H}, Q\right),
\end{aligned}$$

which implies that

$$\begin{aligned}
H(\varepsilon, \mathcal{F}, Q) & \leq \ln \left(\sup_{h \in \mathcal{H}} N\left(\frac{\varepsilon}{2}, \mathcal{G} \circ h, Q\right) \right) + H\left(\frac{\varepsilon}{2L}, \mathcal{H}, Q\right) \\
& = \sup_{h \in \mathcal{H}} H\left(\frac{\varepsilon}{2}, \mathcal{G} \circ h, Q\right) + H\left(\frac{\varepsilon}{2L}, \mathcal{H}, Q\right),
\end{aligned}$$

by continuity. Taking the supremum over all probability measures Q with finite support in \mathcal{X} , we obtain the expected result. \square

Lemma 5.2. *Suppose $\beta \geq 1$. Suppose there exists a constant $C > 0$ and $0 < s \leq \ell/\beta$ such that for all $\varepsilon > 0$ we have $H(\varepsilon, \mathcal{H}) \leq C\varepsilon^{-s}$. Then, there exists a constant $A > 0$ depending only on ℓ, β, L and R such that for all $\varepsilon > 0$ we have*

$$H(\varepsilon, \mathcal{F}) \leq A\varepsilon^{-\ell/\beta}.$$

Proof—According to Lemma 5.1, we only need to prove that there exists a constant $A > 0$ depending only on ℓ, β, L and R such that for all $\varepsilon > 0$ we have

$$\sup_{h \in \mathcal{H}} H(\varepsilon, \mathcal{G} \circ h) \leq A\varepsilon^{-\ell/\beta}.$$

To that aim, fix a probability measure Q with support in \mathcal{X} and $h \in \mathcal{H}$. For all $g \in \mathcal{G}$, let g_h be the restriction of g to $S(h) \cap \mathbf{B}$ and let $\mathcal{G}_h := \{g_h : g \in \mathcal{G}\}$. Then, from the transfer theorem, for all $\varepsilon > 0$ we have

$$H(\varepsilon, \mathcal{G} \circ h, Q) = H(\varepsilon, \mathcal{G}_h, Q \circ h^{-1}),$$

where for any Borel set $A \subset S$: $Q \circ h^{-1}(A) := Q(h^{-1}(A))$. Since $S(h)$ is a vector space of dimension $\ell' \leq \ell$, $S(h) \cap \mathbf{B}$ may be identified to the open Euclidean ball $\mathbf{B}_{\ell'}$ in $\mathbb{R}^{\ell'}$ with radius R and center the origin. Also, $Q \circ h^{-1}$ may be seen as of support in $\mathbf{B}_{\ell'}$ and \mathcal{G}_h as a subset of

$$\mathcal{G}_h := \left\{ g \in \mathcal{C}^{\lfloor \beta \rfloor}(\mathbf{B}_{\ell'}) : \|g\|_{\beta} \leq L \right\},$$

where here $\|\cdot\|_{\beta}$ is defined as in (2.3) but with the Euclidean norm in $\mathbb{R}^{\ell'}$. According to Theorem 9.19 in Kosorok (2008), we know that there exists a constant K' depending only on ℓ', β, R , and L such that for all $\varepsilon > 0$ we have

$$\sup_D H(\varepsilon, \mathcal{G}_{\ell'}, D) \leq K' \varepsilon^{-\ell'/\beta},$$

where the supremum is taken over all probability distributions D in $\mathbb{R}^{\ell'}$. Since $\ell' \leq \ell$ and since $\sup_D H(\varepsilon, \mathcal{G}_{\ell'}, D)$ is equal to 0 for ε sufficiently large, we deduce finally that there exists a constant K depending only on ℓ, β, R , and L such that for all $\varepsilon > 0$ we have

$$\sup_D H(\varepsilon, \mathcal{G}_{\ell'}, D) \leq K \varepsilon^{-\ell/\beta},$$

where the supremum is taken over all probability distributions D in $\mathbb{R}^{\ell'}$. Hence, we deduce that for all $\varepsilon > 0$ we have

$$H(\varepsilon, \mathcal{G} \circ h, Q) \leq K \varepsilon^{-\ell/\beta}.$$

Since the result holds uniformly for $h \in \mathcal{H}$ and for Q with support in \mathcal{X} , we deduce finally that for all $\varepsilon > 0$ we have

$$\sup_{h \in \mathcal{H}} H(\varepsilon, \mathcal{G} \circ h) \leq K \varepsilon^{-\ell/\beta},$$

which completes the proof. \square

Proof of Theorem 3.1

Let $P \in \mathcal{D}$. In the proof, $C > 0$ will denote a constant depending only on ℓ , β , R and L and which value may change from line to line. We denote

$$r_n(x) := \mathbb{E}(Y \mathbf{1}\{|Y| \leq T_n\} | X = x).$$

We have

$$\mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 \leq 2\mathbb{E}(\hat{f}(X) - r_n(X))^2 + 2\mathbb{E}(r_n(X) - \mathbb{E}(Y|X))^2. \quad (5.2)$$

Using Theorem A.1 of the Appendix with $\varepsilon := 1$, $Z := Y \mathbf{1}\{|Y| \leq T_n\}$ and Lemma 5.2, we deduce that there exists a constant $C > 0$ such that

$$\begin{aligned} & \mathbb{E}(\hat{f}(X) - r_n(X))^2 \\ & \leq 2 \inf_{f \in \mathcal{F}} \mathbb{E}(f(X) - r_n(X))^2 + C \left(\frac{(T_n + L)^2}{n} \right)^{2\beta/(2\beta+\ell)} + C \left(\frac{(T_n + L)^2}{n} \right) \\ & \leq 2\mathbb{E}(r_n(X) - \mathbb{E}(Y|X))^2 + C \left(\frac{(T_n + L)^2}{n} \right)^{2\beta/(2\beta+\ell)} + C \left(\frac{(T_n + L)^2}{n} \right), \end{aligned}$$

where in the second inequality we have used the fact that $r \in \mathcal{F}$. Since T_n goes to $+\infty$ and T_n^2/n goes to 0 as n goes to $+\infty$, we deduce that there exists a constant $C > 0$ and an integer n_0 depending only on L , β and ℓ such that $n \geq n_0$ implies that

$$\mathbb{E}(\hat{f}(X) - r_n(X))^2 \leq 2\mathbb{E}(r_n(X) - \mathbb{E}(Y|X))^2 + C \left(\frac{T_n^2}{n} \right)^{2\beta/(2\beta+\ell)}. \quad (5.3)$$

We deduce from (5.2) and (5.3) that there exists a constant $C > 0$ such that for all $n \geq n_0$ we have

$$\mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 \leq 6\mathbb{E}(r_n(X) - \mathbb{E}(Y|X))^2 + C \left(\frac{T_n^2}{n} \right)^{2\beta/(2\beta+\ell)}. \quad (5.4)$$

Using Jensen's inequality and Cauchy-Schwarz's inequality, we have

$$\begin{aligned} \mathbb{E}(r_n(X) - \mathbb{E}(Y|X))^2 &= \mathbb{E}[\mathbb{E}(Y \mathbf{1}\{|Y| \leq T_n\} | X) - \mathbb{E}(Y|X)]^2 \\ &= \mathbb{E}[\mathbb{E}(Y \mathbf{1}\{|Y| > T_n\} | X)]^2 \\ &\leq \mathbb{E}(Y^2 \mathbf{1}\{|Y| > T_n\}) \\ &\leq \sqrt{\mathbb{E}Y^4} \sqrt{\mathbb{P}(|Y| > T_n)}. \end{aligned} \quad (5.5)$$

Then, using the fact that for all $u \in \mathbb{R}$ we have

$$\frac{(\tau u)^4}{4!} \leq e^{\tau|u|},$$

we deduce from the exponential moment condition (2.5) that

$$\mathbb{E}Y^4 \leq \frac{4!}{\tau^4} \mathbb{E}e^{\tau|Y|} \leq \frac{24B}{\tau^4}. \quad (5.6)$$

Using Markov's inequality and the exponential moment condition (2.5), we obtain that

$$\mathbb{P}(|Y| > T_n) = \mathbb{P}\left(e^{\tau|Y|} > e^{\tau T_n}\right) \leq B e^{-\tau T_n}. \quad (5.7)$$

Combining (5.6) and (5.7) we deduce from (5.5) that

$$\mathbb{E}(r_n(X) - \mathbb{E}(Y|X))^2 \leq \frac{\sqrt{24B}}{\tau^2} e^{-\tau T_n/2}. \quad (5.8)$$

Equations (5.4) and (5.8) imply that there exists a constant $C > 0$ such that for all $n \geq n_0$ we have

$$\mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 \leq \frac{6\sqrt{24B}}{\tau^2} e^{-\tau T_n/2} + C \left(\frac{T_n^2}{n}\right)^{2\beta/(2\beta+\ell)}. \quad (5.9)$$

Since the integer n_0 and the constants involved on the right hand side of (5.9) do not depend on $P \in \mathcal{D}$, we deduce that there exists a constant $C > 0$ such that for all $n \geq n_0$ we have

$$\sup_{P \in \mathcal{D}} \mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 \leq \frac{6\sqrt{24B}}{\tau^2} e^{-\tau T_n/2} + C \left(\frac{T_n^2}{n}\right)^{2\beta/(2\beta+\ell)}.$$

Since $\alpha > 2$, the choice of $T_n = (\ln n)^{\alpha/2}$ leads to

$$e^{-\tau T_n/2} \underset{n \rightarrow +\infty}{\ll} \left(\frac{T_n^2}{n}\right)^{2\beta/(2\beta+\ell)} = \mathbf{v}_{n,\ell}^2.$$

It follows that

$$\limsup_{n \rightarrow +\infty} \mathbf{v}_{n,\ell}^{-2} \sup_{P \in \mathcal{D}} \mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 < +\infty.$$

The proof is complete. \square

5.2 Proof of Theorem 3.3

Let \mathcal{D}° be the class of distributions P of (X, Y) where X is of distribution μ , where

$$Y = f(X) + \xi,$$

for $f \in \mathcal{F}$ and where ξ is independent from X and with distribution $\mathcal{N}(0, \sigma^2)$. It may be easily verified that the exponential moment condition (2.5) holds in this context so that we have $\mathcal{D}^\circ \subset \mathcal{D}$. Therefore, it is clear that

$$\inf_{\hat{f}} \sup_{P \in \mathcal{D}^\circ} \mathbb{E} (\hat{f}(X) - \mathbb{E}(Y|X))^2 \leq \inf_{\hat{f}} \sup_{P \in \mathcal{D}} \mathbb{E} (\hat{f}(X) - \mathbb{E}(Y|X))^2.$$

As a result, in order to prove Theorem 3.3 we need only to prove that

$$\liminf_{n \rightarrow +\infty} n^{2\beta/(2\beta+\ell)} \inf_{\hat{f}} \sup_{P \in \mathcal{D}^\circ} \mathbb{E} (\hat{f}(X) - \mathbb{E}(Y|X))^2 > 0. \quad (5.10)$$

According to Theorem 6 in Yang and Barron (1999), inequality (5.10) is satisfied provided there exists a lower bound $N(\varepsilon)$ for the covering number $N(\varepsilon, \mathcal{F}, \mu)$ such that any ε_n satisfying

$$\ln N(\varepsilon_n) = n\varepsilon_n^2,$$

is of order $n^{-\beta/(2\beta+\ell)}$. To obtain such a lower bound, let $h \in \mathcal{H}$ satisfying the conditions of Theorem 3.3. Then, we have $\{g \circ h : g \in \mathcal{G}\} \subset \mathcal{F}$ which implies that for all $\varepsilon > 0$ we have

$$N(\varepsilon, \mathcal{F}, \mu) \geq N(\varepsilon, \{g \circ h : g \in \mathcal{G}\}, \mu). \quad (5.11)$$

The right hand side of inequality (5.11) may be lower bounded as follows. For all $g \in \mathcal{G}$, let g_h be the restriction of g to $S(h) \cap \mathbf{B}$. Then, for all $g, g' \in \mathcal{G}$, we have

$$\begin{aligned} \|g \circ h - g' \circ h\|_\mu^2 &= \int (g \circ h(x) - g' \circ h(x))^2 \mu(dx) \\ &= \int (g_h(u) - g'_h(u))^2 \mu \circ h^{-1}(du) \\ &\geq c \int (g_h(u) - g'_h(u))^2 \lambda_h(du), \end{aligned} \quad (5.12)$$

where in (5.12) we have used (3.9). Hence, for all $\varepsilon > 0$ we have

$$N(\varepsilon, \{g \circ h : g \in \mathcal{G}\}, \mu) \geq N\left(\frac{\varepsilon}{c}, \mathcal{G}_h, \lambda_h\right), \quad (5.13)$$

where we have denoted $\mathcal{G}_h := \{g_h : g \in \mathcal{G}\}$. Now, let us identify $S(h) \cap \mathbf{B}$ to the open Euclidean ball in \mathbb{R}^ℓ with center the origin and radius \mathbb{R}^ℓ and λ_h to the Lebesgue measure in \mathbb{R}^ℓ . Then, according to Corollary 2.4 of chapter 15 in the book by Lorentz et al (1996), there exists a constant $c' > 0$ such that for all $\varepsilon > 0$, we have

$$\ln N\left(\frac{\varepsilon}{c}, \mathcal{G}_h, \lambda_h\right) \geq c' \varepsilon^{-\ell/\beta}. \quad (5.14)$$

Combining (5.11), (5.13) and (5.14) we deduce that

$$\ln N(\varepsilon, \mathcal{F}, \mu) \geq c' \varepsilon^{-\ell/\beta} =: \ln N(\varepsilon).$$

It may be easily verified that the solution ε_n of $c' \varepsilon_n^{-\ell/\beta} = n \varepsilon_n^2$ is given by a constant times $n^{-\beta/(2\beta+\ell)}$, and this concludes the proof. \square

5.3 Proof of Theorem 4.1

Lemma 5.3. *Suppose $\beta \geq 1$ and suppose there exists a constant $C > 0$ and $0 < s \leq p/\beta$ such that for all $\varepsilon > 0$ we have $H(\varepsilon, \mathcal{H}) \leq C\varepsilon^{-s}$. Suppose $a_n = n^u$ and $b_n = n^{-u'}$ with $u > 0$, $u' \geq 0$ and $0 < 2u + u' < 2\beta/(2\beta + p)$. Then for all $\ell \in \{1, \dots, p\}$ and for all $\vartheta > 0$ we have*

$$\lim_{n \rightarrow +\infty} n^\vartheta \sup_{P \in \mathcal{D}} \mathbb{P}\left(|\hat{R}_\ell - R_\ell| \geq b_n\right) = 0.$$

Proof – Let

$$\tilde{R}_\ell := \inf_{f \in \mathcal{F}_\ell} \mathbb{E}(Y^{a_n} - f(X))^2,$$

where $Y^{a_n} := Y \mathbf{1}\{|Y| \leq a_n\}$. We have

$$\begin{aligned} |\hat{R}_\ell - R_\ell| &\leq |\hat{R}_\ell - \tilde{R}_\ell| + |\tilde{R}_\ell - R_\ell| \\ &\leq \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (Y_i^{a_n} - f(X_i))^2 - \mathbb{E}(Y^{a_n} - f(X))^2 \right| \\ &\quad + \sup_{f \in \mathcal{F}} \left| \mathbb{E}(Y^{a_n} - f(X))^2 - \mathbb{E}(Y - f(X))^2 \right|. \end{aligned} \quad (5.15)$$

Using Cauchy-Schwarz's inequality we obtain that for all $f \in \mathcal{F}$, we have

$$\begin{aligned} &\left| \mathbb{E}(Y^{a_n} - f(X))^2 - \mathbb{E}(Y - f(X))^2 \right| \\ &= \left| \mathbb{E}(Y^{a_n} + Y - 2f(X))(Y^{a_n} - Y) \right| \\ &\leq \mathbb{E}[|Y^{a_n} + Y - 2f(X)| |Y^{a_n} - Y|] \\ &\leq \sqrt{\mathbb{E}(Y^{a_n} + Y - 2f(X))^2} \sqrt{\mathbb{E}(Y^{a_n} - Y)^2}. \end{aligned} \quad (5.16)$$

Using the fact that for all $u \in \mathbb{R}$ and for every positive integer k we have

$$\frac{(\tau u)^k}{k!} \leq e^{\tau|u|},$$

we deduce from the exponential moment condition (2.5) that for every positive integer k we have

$$\mathbb{E}Y^k \leq \frac{k!}{\tau^k} \mathbb{E}e^{\tau|Y|} \leq \frac{k!B}{\tau^k}. \quad (5.17)$$

Then, using Minkowski's inequality and the fact that all functions $f \in \mathcal{F}$ are bounded by L , we deduce that

$$\begin{aligned} \sqrt{\mathbb{E}(Y^{a_n} + Y - 2f(X))^2} &\leq 2\sqrt{\mathbb{E}Y^2} + 2\sqrt{\mathbb{E}f(X)^2} \\ &\leq 2\sqrt{\frac{2!}{\tau^2} \mathbb{E}e^{\tau|Y|}} + 2L \\ &\leq 2\sqrt{\frac{2B}{\tau^2}} + 2L. \end{aligned} \quad (5.18)$$

Using Markov's inequality and the exponential moment condition (2.5) we obtain that

$$\mathbb{P}(|Y| > a_n) = \mathbb{P}\left(e^{\tau|Y|} > e^{\tau a_n}\right) \leq B e^{-\tau a_n}. \quad (5.19)$$

Then, using Cauchy-Schwarz's inequality and equations (5.17) and (5.19) we deduce that

$$\begin{aligned} \mathbb{E}(Y^{a_n} - Y)^2 &= \mathbb{E}(Y^2 \mathbf{1}\{|Y| > a_n\}) \\ &\leq \sqrt{\mathbb{E}Y^4} \sqrt{\mathbb{P}(|Y| > a_n)} \\ &\leq \sqrt{\frac{4!}{\tau^4} \mathbb{E}e^{\tau|Y|}} \sqrt{\mathbb{P}(e^{\tau|Y|} > e^{\tau a_n})} \\ &\leq B \sqrt{\frac{24}{\tau^4}} e^{-\tau a_n/2}. \end{aligned} \quad (5.20)$$

Hence, combining (5.16), (5.18) and (5.20) we deduce that

$$\begin{aligned} \sup_{f \in \mathcal{F}} \left| \mathbb{E}(Y^{a_n} - f(X))^2 - \mathbb{E}(Y - f(X))^2 \right| &\leq \sqrt{\left(\frac{2\sqrt{2B}}{\tau} + 2L\right)} \sqrt{\left(\frac{B\sqrt{24}}{\tau^2}\right)} e^{-\tau a_n/4} \\ &=: U e^{-\tau a_n/4}. \end{aligned}$$

Therefore, denoting by $\kappa_n := b_n - Ue^{-\tau a_n/4}$ and

$$\bar{\mathcal{F}}^{a_n} := \left\{ (x, y) \mapsto (y \mathbf{1}\{|y| \leq a_n\} - f(x))^2 : f \in \mathcal{F} \right\},$$

we deduce from (5.15) and Theorem 9.1 in the book by Györfi *et al.* (2002) that there exists a universal constant $C > 0$ such that

$$\begin{aligned} \mathbb{P}(|\hat{R}_\ell - R_\ell| \geq b_n) &\leq \mathbb{P}\left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (Y_i^{a_n} - f(X_i))^2 - \mathbb{E}(Y^{a_n} - f(X))^2 \right| \geq \kappa_n\right) \\ &\leq C \mathbb{E} \left[N_1\left(\frac{\kappa_n}{C}, \bar{\mathcal{F}}^{a_n}, P_n\right) \right] e^{-\frac{n\kappa_n^2}{C(a_n+L)^4}}, \end{aligned} \quad (5.21)$$

where $P_n := n^{-1} \sum_{i=1}^n \delta_{(X_i, Y_i)}$ denotes the empirical distribution associated with the sample $(X_1, Y_1), \dots, (X_n, Y_n)$ and where $N_1(\varepsilon, \mathcal{C}, Q)$ denotes the minimal number of metric balls of radius ε in $\mathbb{L}^1(Q)$ that are needed to cover \mathcal{C} . For all $f, f' \in \mathcal{F}$ we have

$$\begin{aligned} &\frac{1}{n} \sum_{i=1}^n \left| (Y_i^{a_n} - f(X_i))^2 - (Y_i^{a_n} - f'(X_i))^2 \right| \\ &= \frac{1}{n} \sum_{i=1}^n |2Y_i^{a_n} - f(X_i) - f'(X_i)| |f(X_i) - f'(X_i)| \\ &\leq \frac{2(a_n + L)}{n} \sum_{i=1}^n |f(X_i) - f'(X_i)| \\ &\leq 2(a_n + L) \left\{ \frac{1}{n} \sum_{i=1}^n (f(X_i) - f'(X_i))^2 \right\}^{1/2}. \end{aligned}$$

Therefore, we obtain that

$$N_1\left(\frac{\kappa_n}{C}, \bar{\mathcal{F}}^{a_n}, P_n\right) \leq N\left(\frac{\kappa_n}{2C(a_n+L)}, \mathcal{F}, \mu_n\right),$$

where $\mu_n := n^{-1} \sum_{i=1}^n \delta_{X_i}$. Hence, we deduce from (5.21) and Lemma 5.2 that there exists a universal constant $C > 0$ such that

$$\mathbb{P}\left(|\hat{R}_\ell - R_\ell| \geq b_n\right) \leq C \exp\left(\left(\frac{C(a_n+L)}{\kappa_n}\right)^{p/\beta} - \frac{n\kappa_n^2}{C(a_n+L)^4}\right). \quad (5.22)$$

Recalling that $a_n = n^u$, that $b_n = n^{-u'}$ and that $\kappa_n = b_n - Ue^{-\tau a_n/4}$, it follows that there exists a universal constant $C > 0$ such that

$$\sup_{P \in \mathcal{D}} \mathbb{P}\left(|\hat{R}_\ell - R_\ell| \geq b_n\right) \leq C \exp\left(n^{(u+u')p/\beta} - n^{1-2u'-4u}\right). \quad (5.23)$$

Consequently, the term on the left hand side of (5.23) converges exponentially fast to 0 as n goes to $+\infty$ provided

$$(1 - 2u' - 4u) > (u + u') \frac{p}{\beta},$$

which is implied by

$$0 < 2u + u' < \frac{2\beta}{2\beta + p}.$$

The proof is complete. \square

Proof of Theorem 4.1

Let $P \in \mathcal{D}$. Since the function $\ell \mapsto \hat{R}_\ell$ is non increasing, for all integer $q \in \{1, \dots, p\}$ and all $n \geq 1$ we have

$$\min \left\{ \ell = 1, \dots, p : \hat{R}_\ell - \hat{R}_p \leq b_n \right\} \leq q \Leftrightarrow \hat{R}_q - \hat{R}_p \leq b_n. \quad (5.24)$$

Therefore, using the fact that $R_d = R_p$, we obtain that

$$\begin{aligned} \mathbb{P}(\hat{d} > d) &= \mathbb{P}\left(\min \left\{ \ell = 1, \dots, p : \hat{R}_\ell - \hat{R}_p \leq b_n \right\} > d\right) \\ &= \mathbb{P}\left(\hat{R}_d - \hat{R}_p > b_n\right) \\ &= \mathbb{P}\left((\hat{R}_d - R_d) + (R_p - \hat{R}_p) > b_n\right) \\ &\leq \mathbb{P}\left(|\hat{R}_d - R_d| \geq \frac{b_n}{2}\right) + \mathbb{P}\left(|\hat{R}_p - R_p| \geq \frac{b_n}{2}\right). \end{aligned}$$

Using Lemma 5.3, we deduce that for all $\vartheta > 0$ we have

$$\lim_{n \rightarrow +\infty} n^\vartheta \sup_{P \in \mathcal{D}} \mathbb{P}\left(|\hat{R}_d - R_d| \geq \frac{b_n}{2}\right) = 0,$$

and

$$\lim_{n \rightarrow +\infty} n^\vartheta \sup_{P \in \mathcal{D}} \mathbb{P}\left(|\hat{R}_p - R_p| \geq \frac{b_n}{2}\right) = 0,$$

which gives

$$\lim_{n \rightarrow +\infty} n^\vartheta \sup_{P \in \mathcal{D}} \mathbb{P}(\hat{d} > d) = 0. \quad (5.25)$$

Now let $P \in \mathcal{D}(\underline{\delta})$ and assume $d \geq 2$. Using the fact that $R_{d-1} - R_p = \Delta \geq \underline{\delta}$ and provided n is large enough to have $\underline{\delta} - b_n \geq b_n$, we obtain that

$$\begin{aligned}
\mathbb{P}(\hat{d} < d) &= \mathbb{P}\left(\min\{\ell = 1, \dots, p : \hat{R}_\ell - \hat{R}_p \leq b_n\} \leq d-1\right) \\
&= \mathbb{P}\left(\hat{R}_{d-1} - \hat{R}_p \leq b_n\right) \\
&= \mathbb{P}\left((\hat{R}_{d-1} - R_{d-1}) + \Delta + (R_p - \hat{R}_p) \leq b_n\right) \\
&\leq \mathbb{P}\left((R_{d-1} - \hat{R}_{d-1}) + (\hat{R}_p - R_p) \geq \underline{\delta} - b_n\right) \\
&\leq \mathbb{P}\left((R_{d-1} - \hat{R}_{d-1}) + (\hat{R}_p - R_p) \geq b_n\right) \\
&\leq \mathbb{P}\left(|\hat{R}_{d-1} - R_{d-1}| \geq \frac{b_n}{2}\right) + \mathbb{P}\left(|\hat{R}_p - R_p| \geq \frac{b_n}{2}\right).
\end{aligned}$$

From the same argument as in the beginning of the proof, we have

$$\lim_{n \rightarrow +\infty} n^\vartheta \sup_{P \in \mathcal{D}(\underline{\delta})} \mathbb{P}(\hat{d} < d) = 0,$$

which concludes the proof. \square

5.4 Proof of Theorem 4.2

Since we have

$$\inf\{\Delta : P \in \mathcal{D}_0\} = 0,$$

for all $\varepsilon > 0$ there exists $P(\varepsilon) \in \mathcal{D}_0$ such that

$$\Delta_{P(\varepsilon)} \leq \varepsilon.$$

For all $n \geq 1$ let

$$Q_n := P(b_n/2).$$

Now assume that the sample $(X_1, Y_1), \dots, (X_n, Y_n)$ is drawn from Q_n and let \hat{d} be defined by

$$\hat{d} := \min\{\ell : \hat{R}_\ell \leq \hat{R}_p + b_n\}.$$

Then, we obtain as in the proof of Theorem 4.1 that

$$\begin{aligned}
\mathbb{P}(\hat{d} < d) &= \mathbb{P}\left(\min\{\ell = 1, \dots, p : \hat{R}_\ell - \hat{R}_p \leq b_n\} \leq d-1\right) \\
&= \mathbb{P}\left(\hat{R}_{d-1} - \hat{R}_p \leq b_n\right) \\
&= \mathbb{P}\left((\hat{R}_{d-1} - R_{d-1}) + \Delta_{Q_n} + (R_p - \hat{R}_p) \leq b_n\right) \\
&\geq \mathbb{P}\left((\hat{R}_{d-1} - R_{d-1}) + \frac{b_n}{2} + (R_p - \hat{R}_p) \leq b_n\right) \\
&= \mathbb{P}\left((\hat{R}_{d-1} - R_{d-1}) + (R_p - \hat{R}_p) \leq \frac{b_n}{2}\right) \\
&\geq 1 - \mathbb{P}\left((\hat{R}_{d-1} - R_{d-1}) + (R_p - \hat{R}_p) \geq \frac{b_n}{2}\right) \\
&\geq 1 - \mathbb{P}\left(|\hat{R}_{d-1} - R_{d-1}| \geq \frac{b_n}{4}\right) - \mathbb{P}\left(|\hat{R}_p - R_p| \geq \frac{b_n}{4}\right).
\end{aligned}$$

According to Lemma 5.3, we know that for all $\vartheta > 0$ we have

$$\lim_{n \rightarrow +\infty} n^\vartheta \mathbb{P}\left(|\hat{R}_{d-1} - R_{d-1}| \geq \frac{b_n}{4}\right) = 0 \quad \text{and} \quad \lim_{n \rightarrow +\infty} n^\vartheta \mathbb{P}\left(|\hat{R}_p - R_p| \geq \frac{b_n}{4}\right) = 0.$$

As a result, there exists an integer n_0 such that for all $n \geq n_0$ we have

$$\mathbb{P}\left(|\hat{R}_{d-1} - R_{d-1}| \geq \frac{b_n}{4}\right) \leq \frac{1}{2n^\vartheta} \quad \text{and} \quad \mathbb{P}\left(|\hat{R}_p - R_p| \geq \frac{b_n}{4}\right) \leq \frac{1}{2n^\vartheta}.$$

Therefore, for all $n \geq n_0$, we have

$$\sup_{P \in \mathcal{D}_0} \mathbb{P}(\hat{d} < d) \geq \mathbb{P}_{(X,Y) \sim Q_n}(\hat{d} < d) \geq 1 - \frac{1}{n^\vartheta},$$

which concludes the proof. \square

5.5 Proof of Theorem 4.3

Let $\underline{\delta} > 0$ and $P \in \mathcal{D}(\underline{\delta})$. We have

$$\begin{aligned}
\mathbb{E}(\hat{f}(X) - \mathbb{E}(X|Y))^2 &= \mathbb{E}\left[\mathbf{1}\{\hat{d} \neq d\} (\hat{f}(X) - \mathbb{E}(X|Y))^2\right] \\
&\quad + \mathbb{E}\left[\mathbf{1}\{\hat{d} = d\} (\hat{f}(X) - \mathbb{E}(X|Y))^2\right]. \quad (5.26)
\end{aligned}$$

Since both \hat{f} and r belong to \mathcal{F} , they are bounded by L and therefore we have

$$\mathbb{E}\left[\mathbf{1}\{\hat{d} \neq d\} (\hat{f}(X) - \mathbb{E}(X|Y))^2\right] \leq 4L^2 \mathbb{P}(\hat{d} \neq d). \quad (5.27)$$

Then, we observe that

$$\begin{aligned} \mathbb{E}\left[\mathbf{1}\{\hat{d} = d\} (\hat{f}(X) - \mathbb{E}(X|Y))^2\right] &= \mathbb{E}\left[\mathbf{1}\{\hat{d} = d\} (\hat{f}_d(X) - \mathbb{E}(Y|X))^2\right] \\ &\leq \mathbb{E}(\hat{f}_d(X) - \mathbb{E}(Y|X))^2. \end{aligned} \quad (5.28)$$

Hence, we deduce from (5.26), (5.27) and (5.28) that

$$\begin{aligned} \mathbf{v}_{n,d}^{-2} \mathbb{E}(\hat{f}(X) - \mathbb{E}(X|Y))^2 &\leq 4L^2 \mathbf{v}_{n,d}^{-2} \mathbb{P}(\hat{d} \neq d) \\ &\quad + \mathbf{v}_{n,d}^{-2} \mathbb{E}(\hat{f}_d(X) - \mathbb{E}(Y|X))^2. \end{aligned} \quad (5.29)$$

According to the proof of Theorem 3.1, there exists a constant $C > 0$ depending only on τ, B, β, R and L such that

$$\mathbf{v}_{n,d}^{-2} \mathbb{E}(\hat{f}_d(X) - \mathbb{E}(Y|X))^2 \leq C. \quad (5.30)$$

In particular, since this constant C does not depend on $P \in \mathcal{D}$ we obtain that

$$\sup_{P \in \mathcal{D}(\underline{\delta})} \mathbf{v}_{n,d}^{-2} \mathbb{E}(\hat{f}_d(X) - \mathbb{E}(Y|X))^2 \leq C. \quad (5.31)$$

Then, we deduce from Theorem 4.1 that

$$\sup_{P \in \mathcal{D}(\underline{\delta})} \mathbf{v}_{n,d}^{-2} \mathbb{P}(\hat{d} \neq d) \leq \mathbf{v}_{n,1}^{-2} \sup_{P \in \mathcal{D}(\underline{\delta})} \mathbb{P}(\hat{d} \neq d) \xrightarrow{n \rightarrow +\infty} 0. \quad (5.32)$$

Finally, it follows from (5.29), (5.31) and (5.32) that there exists a constant depending only on τ, B, β, R and L such that

$$\limsup_{n \rightarrow +\infty} \sup_{P \in \mathcal{D}(\underline{\delta})} \mathbf{v}_{n,d}^{-2} \mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 \leq C,$$

which concludes the proof. \square

A Performance of least-squares estimates

Let (X, Z) be an $\mathcal{X} \times \mathbb{R}$ -valued random variable of distribution P where \mathcal{X} denotes a metric space. We fix $T > 0$ and assume that $|Z| \leq T$. Let μ be the distribution of X . For $x \in \mathcal{X}$ we set

$$f_*(x) := \mathbb{E}(Z|X = x). \quad (\text{A.1})$$

Let $L > 0$ be fixed and \mathcal{F} be a class of functions $f : \mathcal{X} \rightarrow \mathbb{R}$ such that

$$\sup_{x \in \mathcal{X}} |f(x)| \leq L. \quad (\text{A.2})$$

Let $(X_1, Z_1), \dots, (X_n, Z_n)$ be a sample of n independent and identically distributed random variables with distribution P . The least-squares estimate f_n of f_* is defined as any element in \mathcal{F} satisfying

$$f_n \in \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (Z_i - f(X_i))^2.$$

Implicitly, it is assumed that such an element exists. Following the lines devised in Koltchinskii (2006), we establish the following result. The proof may be found in the supplementary material Paris (2012).

Theorem A.1. *Let $0 < s < 2$ and assume there exists $A > 0$ such that for all $u > 0$ we have $H(u, \mathcal{F}) \leq Au^{-s}$. Then, there exists a constant C depending only on s, A and L such that for all $\varepsilon \in (0, 1]$ we have*

$$\mathbb{E} \|f_n - f_*\|_{\mu}^2 \leq (1 + \varepsilon) \inf_{f \in \mathcal{F}} \|f - f_*\|_{\mu}^2 + C\varepsilon^{-\frac{2-s}{2+s}} \left(\frac{b}{n}\right)^{\frac{2}{2+s}} + \frac{C}{\varepsilon} \left(\frac{b}{n}\right),$$

where $b := (T + L)^2$.

B Reduced dimension d and parameter Δ

In this appendix, we prove equations (4.5) and (4.8). First, observe that since \mathcal{F}_ℓ is compact in $\mathbb{L}^2(\mu)$ and since $r \in \mathcal{F}$, we have

$$\begin{aligned}
 r \in \mathcal{F}_\ell &\Leftrightarrow \inf_{f \in \mathcal{F}_\ell} \mathbb{E}(f(X) - r(X))^2 = 0 \\
 &\Leftrightarrow \inf_{f \in \mathcal{F}_\ell} \mathbb{E}(Y - f(X))^2 - \mathbb{E}(Y - r(X))^2 = 0 \\
 &\Leftrightarrow \inf_{f \in \mathcal{F}_\ell} \mathbb{E}(Y - f(X))^2 - \inf_{f \in \mathcal{F}} \mathbb{E}(Y - f(X))^2 = 0 \\
 &\Leftrightarrow R_\ell = R_p.
 \end{aligned}$$

Therefore, since the function $\ell \in \{1, \dots, p\} \mapsto R_\ell$ is non-increasing, we deduce that

$$\begin{aligned}
 d &:= \min \left\{ \ell : r \in \mathcal{F}_\ell \right\} \\
 &= \min \left\{ \ell : R_\ell = R_p \right\},
 \end{aligned}$$

which proves equation (4.5). Using (4.5) and the fact that $r \in \mathcal{F}$ we obtain that

$$\begin{aligned}
 \Delta &= \min \left\{ R_\ell - R_p : R_\ell > R_p \right\} \\
 &= R_{d-1} - R_p \\
 &= \inf_{f \in \mathcal{F}_{d-1}} \mathbb{E}(Y - f(X))^2 - \inf_{f \in \mathcal{F}} \mathbb{E}(Y - f(X))^2 \\
 &= \inf_{f \in \mathcal{F}_{d-1}} \mathbb{E}(Y - f(X))^2 - \mathbb{E}(Y - r(X))^2 \\
 &= \inf_{f \in \mathcal{F}_{d-1}} \mathbb{E}(f(X) - r(X))^2 \\
 &= \inf_{f \in \mathcal{F}_{d-1}} \|f - r\|_\mu^2,
 \end{aligned}$$

which proves (4.8).

Acknowledgments – The author is indebted to Benoit Cadre and Bruno Pelletier for their time and help. Nicolas Klutchnikoff is also to be thanked for several fruitful discussions and advices concerning the redaction of the manuscript.

References

- Cadre, B. and Dong, Q. (2010). Dimension reduction in regression estimation with nearest neighbor. *Electronic Journal of Statistics*. Vol. 4, pp. 436-460.
- Cook, R.D. and Li, B. (2002). Dimension reduction for conditional mean in regression. *The Annals of Statistics*. Vol. 30, pp. 455-474.
- Cucker, F. and Smale, S. (2001). On the mathematical foundations of learning. *Bulletin of the American Mathematical Society*. Vol. 39, pp. 1-49.
- Györfi, L., Kohler, M., Krzyzak, A. and Walk, H. (2002). *A Distribution-Free Theory of Nonparametric Regression*. Springer-Verlag, New-York NY.
- Ibragimov, I.A. and Khasminskii, R.Z. (1981). *Statistical Estimation: Asymptotic Theory*. Springer-Verlag, New-York NY.
- Kohler, M., Krzyzak, A. and Walk, H. (2009). Optimal global rates of convergence in nonparametric regression with unbounded data. *Journal of Statistical Planning and Inference*. Vol. 123, pp. 1286-1296.
- Koltchinskii, V. (2006). Local Rademacher complexities and oracle inequalities in risk minimization. *The Annals of Statistics*. Vol. 34, pp. 2593-2656.
- Koltchinskii, V. (2011). *Oracle Inequalities in Empirical Risk Minimization and Sparse Recovery Problems*. Ecole d'été de Probabilités de Saint-Flour 2008. Lectures Notes in Mathematics 2033.
- Kosorok, M.R. (2008). *Introduction to Empirical Processes and Semiparametric Inference*. Springer-Verlag, New-York NY.
- Li, K.C. (1991). Sliced inverse regression for dimension reduction (with discussion). *Journal of the American Statistical Association*. Vol. 86, pp. 316-342.
- Lorentz, G.G., Golitschek, M.v. and Makovoz, Y. (1996). *Constructive Approximation: Advanced Problems*. Springer-Verlag, New-York NY.
- Paris, Q. (2012). Supplement to "Minimax adaptative dimension reduction for regression".
- Yang, Y. and Barron, A. (1999). Information-theoretic determination of minimax rates of convergence. *The Annals of Statistics*. Vol. 27, pp. 1564-1599.