



# Minimax adaptive dimension reduction for regression

Quentin Paris

IRMAR, ENS Cachan Bretagne, CNRS, UEB, Campus de Ker Lann, Avenue Robert Schuman, 35170 Bruz, France

## ARTICLE INFO

### Article history:

Received 26 December 2012

Available online 4 April 2014

### AMS 2000 subject classifications:

62H12

62G08

### Keywords:

Regression estimation

Dimension reduction

Minimax rates of convergence

Empirical risk minimization

Metric entropy

## ABSTRACT

In this paper, we address the problem of regression estimation in the context of a  $p$ -dimensional predictor when  $p$  is large. We propose a general model in which the regression function is a composite function. Our model consists in a nonlinear extension of the usual sufficient dimension reduction setting. The strategy followed for estimating the regression function is based on the estimation of a new parameter, called the reduced dimension. We adopt a minimax point of view and provide both lower and upper bounds for the optimal rates of convergence for the estimation of the regression function in the context of our model. We prove that our estimate adapts, in the minimax sense, to the unknown value  $d$  of the reduced dimension and achieves therefore fast rates of convergence when  $d \ll p$ .

© 2014 Elsevier Inc. All rights reserved.

## Contents

1. Introduction.....	187
1.1. The curse of dimensionality in regression .....	187
1.2. A general model for dimension reduction in regression .....	187
1.3. Organization of the paper .....	188
2. Model and statistical methodology.....	188
2.1. The model.....	188
2.2. The reduced dimension .....	188
2.3. Estimation of the regression function .....	189
3. Results.....	190
3.1. Performance of the estimates $\hat{r}_\ell$ .....	190
3.2. Performance of $\hat{d}$ .....	191
3.3. Dimension adaptivity of $\hat{r}$ .....	192
4. Proofs.....	193
4.1. Proof of Theorem 3.1 .....	193
4.2. Proof of Theorem 3.2 .....	195
4.3. Proof of Theorem 3.3 .....	196
4.4. Proof of Theorem 3.4 .....	199
4.5. Proof of Theorem 3.5 .....	200
Acknowledgments .....	200
Appendix A. ....	200
A.1. Reduced dimension $d$ and parameter $\Delta$ .....	200
A.2. Performance of least-squares estimates .....	201
Appendix B. Supplementary material.....	201
References.....	201

E-mail address: [quentin.paris@bretagne.ens-cachan.fr](mailto:quentin.paris@bretagne.ens-cachan.fr).

<http://dx.doi.org/10.1016/j.jmva.2014.03.008>

0047-259X/© 2014 Elsevier Inc. All rights reserved.

## 1. Introduction

### 1.1. The curse of dimensionality in regression

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 the generic random couple  $(X, Y)$ .

A major issue in regression, known as the curse of dimensionality, is basically that the rates of convergence of estimates of the regression function are slow when the dimension  $p$  of the predictor variable  $X$  is high. For instance, if  $r$  is assumed to be  $\beta$ -Hölder and if  $\hat{r}$  refers to any classical estimate (say a kernel, a nearest-neighbors or a least-squares estimate), the mean squared error  $\mathbb{E}(\hat{r}(X) - r(X))^2$  of  $\hat{r}$  converges to 0 at the rate  $n^{-2\beta/(2\beta+p)}$ , which gets slower as  $p$  increases. To get a deeper understanding of the problem, one may refer to the minimax point of view. First, we recall the definition of optimal rates of convergence in the minimax sense. Given a set  $\mathcal{D}$  of distributions  $P$  of the random couple  $(X, Y)$ ,  $v_n$  is said to be an optimal rate of convergence in the minimax sense for  $\mathcal{D}$  if it is a lower minimax rate, i.e.,

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

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

$$\limsup_{n \rightarrow +\infty} v_n^{-2} \sup_{P \in \mathcal{D}} \mathbb{E}(\hat{r}(X) - r(X))^2 < +\infty.$$

Then, in a word, when  $\mathcal{D}$  is taken as the set of all distributions  $P$  of the random couple  $(X, Y)$  for which  $r$  is  $\beta$ -Hölder, the optimal rate of convergence for  $\mathcal{D}$  is  $v_n = n^{-\beta/(2\beta+p)}$  (for more details on optimal rates of convergence, we refer the reader to [22,12,14,23]). Accordingly, there is no hope of constructing an estimate which converges at a faster rate under the only general assumption that  $r$  is regular. Hence, the only alternative to obtain faster rates is to exploit additional information on the regression function.

### 1.2. A general model for dimension reduction in regression

In practice, when such additional information is available, it is often encoded in regression models as so called structural assumptions on the regression function. Statistical procedures based on such models are usually referred to as dimension reduction techniques. In the recent years, much attention has been paid to dimension reduction techniques due to the increasing complexity of the data considered in applications. Among popular models for dimension reduction in regression, one can mention for example the single index model (see, e.g., [1], and the references therein), the additive regression model or the projection pursuit model (see, e.g., Chapter 22 in [12]). Another important dimension reduction framework is called sufficient dimension reduction. In this framework, one assumes that

$$\mathbb{E}(Y|X) = \mathbb{E}(Y|AX) \quad \text{and} \quad \mathbb{E}(Y|AX = \cdot) \in \mathcal{G}, \quad (1.2)$$

are satisfied for a matrix  $A \in \mathcal{M}_p(\mathbb{R})$  of rank smaller than  $p$ , and a class  $\mathcal{G}$  of regular functions (see, e.g., [13,18,3], and the references therein). The motivation for studying such a model is that, provided the matrix  $A$  may be estimated, the predictor variable  $X$  may be replaced by  $AX$  which takes its values in a lower dimensional space. Many methods have been introduced in the literature to estimate  $A$  among which we mention average derivative estimation (ADE) [13], sliced inverse regression (SIR) [18], principal Hessian directions (PHD) [19], sliced average variance estimation (SAVE) [7], kernel dimension reduction (KSIR) [10] and, more recently, the optimal transformation procedure [9]. Discussions, improvements and other relevant papers on that topic can be found in [5,11,26,6,29], and in the references therein. In the last years, little attention has been paid to measuring the impact of the estimation of  $A$  in terms of the estimation of  $r$ . Recently, Cadre and Dong [2] have used these methods to show that, in the context of model (1.2), one could indeed construct an estimate  $\hat{r}$  of the regression function such that

$$\mathbb{E}(\hat{r}(X) - r(X))^2 = O(n^{-2/(2+\text{rank}(A))}),$$

when  $\mathcal{G}$  is taken as a class of Lipschitz functions.

In the present article, we tackle the problem of dimension reduction for regression by studying a model which consists in a nonlinear extension of (1.2) and which is described as follows.

*Our model*—For a given class  $\mathcal{H}$  of functions  $h : \mathcal{X} \rightarrow \mathbb{R}^p$  and a given class  $\mathcal{G}$  of regular functions  $g : \mathbb{R}^p \rightarrow \mathbb{R}$ , we assume that the two conditions

$$(i) \mathbb{E}(Y|X) = \mathbb{E}(Y|h(X)) \quad \text{and} \quad (ii) \mathbb{E}(Y|h(X) = \cdot) \in \mathcal{G}, \quad (1.3)$$

Download English Version:

<https://daneshyari.com/en/article/1145558>

Download Persian Version:

<https://daneshyari.com/article/1145558>

[Daneshyari.com](https://daneshyari.com)