

COBRA: A Nonlinear Aggregation Strategy

Gérard Biau

Université Pierre et Marie Curie¹ & Ecole Normale Supérieure², France
gerard.biau@upmc.fr

Aurélie Fischer

Université Paris Diderot, France
aurelie.fischer@univ-paris-diderot.fr

Benjamin Guedj³

Université Pierre et Marie Curie, France
benjamin.guedj@upmc.fr

James D. Malley

National Institute of Health, U.S.A.
jmalley@mail.nih.gov

Abstract

A new method for combining several initial estimators of the regression function is introduced. Instead of building a linear or convex optimized combination over a collection of basic estimators r_1, \dots, r_M , we use them as a collective indicator of the proximity between the training data and a test observation. This local distance approach is model-free and very fast. More specifically, the resulting collective estimator is shown to perform asymptotically at least as well in the L^2 sense as the best basic estimator in the collective. Moreover, it does so without having to declare which might be the best basic estimator for the given data set. A companion R package called **COBRA** (standing for COmBined Regression Alternative) is presented (downloadable on <http://cran.r-project.org/web/packages/COBRA/index.html>). Substantial numerical evidence is provided on both synthetic and real data sets to assess the excellent performance and velocity of our method in a large variety of prediction problems.

Index terms — Regression estimation, aggregation, nonlinearity, consistency, prediction.

2010 Mathematics Subject Classification: 62G05, 62G20.

¹Research partially supported by the French National Research Agency (grant ANR-09-BLAN-0051-02 “CLARA”) and by the Institut universitaire de France.

²Research carried out within the INRIA project “CLASSIC” hosted by Ecole Normale Supérieure and CNRS.

³Corresponding author.

1 Introduction

Recent years have witnessed a growing interest in aggregated statistical procedures, supported by a considerable research and thorough empirical evidence. Indeed, the increasing number of available estimation and prediction methods (hereafter also denoted as machines) in a wide range of modern statistical problems naturally suggests using some efficient strategy for combining procedures and estimators. If the combined strategy is known to be optimal in some sense and relatively free of assumptions that are hard to evaluate, then such a model-free strategy is a valuable research tool.

In this regard, numerous contributions have enriched the aggregation literature with various approaches, such as model selection aggregation (select the optimal single estimator from a list), convex aggregation (searching for the optimal convex combination of given estimators, such as exponentially weighted aggregates) and linear aggregation (selecting the optimal linear combination).

Model selection, linear-type aggregation strategies and related problems have been studied by [Catoni \(2004\)](#), [Juditsky and Nemirovski \(2000\)](#), [Nemirovski \(2000\)](#), [Yang \(2000, 2001, 2004\)](#), [Györfi et al. \(2002\)](#), and [Wegkamp \(2003\)](#). Minimax results have been derived by [Nemirovski \(2000\)](#) and [Tsybakov \(2003\)](#), leading to the notion of optimal rates of aggregation. Similar results can be found in [Bunea et al. \(2007a\)](#). Further upper bounds for the risk in model selection and convex aggregation have been established for instance by [Audibert \(2004\)](#), [Birgé \(2006\)](#), and [Dalalyan and Tsybakov \(2008\)](#). An interesting feature is that such aggregation problems may be treated within the scope of L^1 -penalized least squares, as performed in [Bunea et al. \(2006, 2007a,b\)](#). This kind of framework is also considered by [van de Geer \(2008\)](#) and [Koltchinskii \(2009\)](#), with the L^2 loss replaced by another convex loss. In the aggregation literature, let us also mention the work of [Juditsky et al. \(2005\)](#), [Bunea and Nobel \(2008\)](#), and [Baraud et al. \(2013\)](#). More recently, specific models such as single-index in [Alquier and Biau \(2013\)](#) and additive models in [Guedj and Alquier \(2013\)](#) have been studied in the context of aggregation under a sparsity assumption.

The present article investigates a distinctly different point of view, motivated by the sense that nonlinear, data-dependent techniques are a source of analytic flexibility and might improve over current aggregation procedures. In this regard, consider the following example regarding classification problem: If the ensemble of machines happens to include a strong one, lurking but unnamed in the collection of which many might be very weak machines, it might make sense to consider a more sophisticated method than the previously cited ones for pooling the information across the machines. Choosing to set aside some of the machines, on some data-dependent criteria, seems only weakly motivated, since the performance of the collective, retaining those suspect

machines, might be quite good on a nearby data set. Similarly, searching for some phantom strong machine in the collective could also be ruinous when presented with new and different data.

Instead of choosing either of these options—selecting out weak performers, searching for a hidden, universally strong performer—we propose an original nonlinear method for combining the outcomes over some list of plausibly good procedures. We call this combined scheme a regression collective over the given basic machines. More specifically, we consider the problem of building a new estimator by combining M estimators of the regression function, thereby exploiting an idea proposed in the context of classification by [Mojirsheibani \(1999\)](#). Given a set of preliminary estimators r_1, \dots, r_M , the idea behind this aggregation method is a “unanimity” concept, in that it is based on the values predicted by r_1, \dots, r_M for the data and for a new observation \mathbf{x} . In a nutshell, a data point is considered to be “close” to \mathbf{x} , and consequently, reliable for contributing to the estimation of this new observation, if all estimators predict values which are close to each other for \mathbf{x} and this data item, *i.e.*, not more distant than a prespecified threshold ε . The predicted value corresponding to this query point \mathbf{x} is then set to the average of the responses of the selected observations. More precisely, the average is over the original outcome values of the selected observations, and *not* over the estimates provided by the several machines for these observations.

To make the concept clear, consider the following toy example illustrated by [Figure 1](#). Assume we are given the observations plotted in circles, and the values predicted by two known machines f_1 and f_2 (triangles pointing up and down, respectively). The goal is to predict the response for the new point \mathbf{x} (along the dotted line). Set a threshold ε , the black solid circles are the data points (\mathbf{x}_i, y_i) within the two dotted intervals, *i.e.*, such that for $m = 1, 2$, $|f_m(\mathbf{x}_i) - f_m(\mathbf{x}_0)| \leq \varepsilon$. Averaging the corresponding y_i ’s yields the prediction for \mathbf{x} (diamond).

We stress that the central and original idea behind our approach is that the resulting regression predictor is a nonlinear, data-dependent function of the basic predictors r_1, \dots, r_M . To the best of our knowledge there exists no formalized procedure in the machine learning and aggregation literature that operates as does ours.

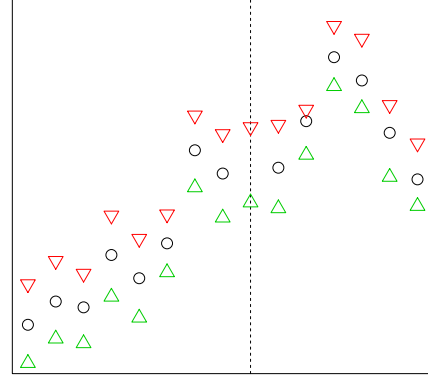
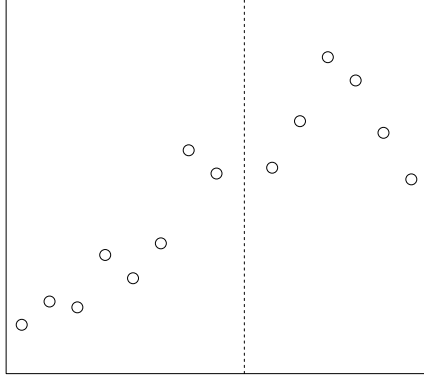
Along with this paper, we release the software **COBRA** ([Guedj, 2013](#)) which implements the method as an additional package to the statistical software R (see [R Core Team, 2013](#)). **COBRA** is freely downloadable on the CRAN web-site⁴. As detailed in [Section 3](#), we undertook a lengthy series of numerical experiments, over which **COBRA** proved extremely successful. These stunning results lead us to believe that regression collectives can provide valuable insights on a wide range of prediction problems. Further, these same results

⁴<http://cran.r-project.org/web/packages/COBRA/index.html>

Figure 1: A toy example: Nonlinear aggregation of two primal estimators.

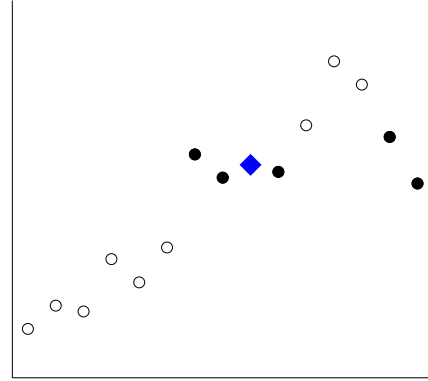
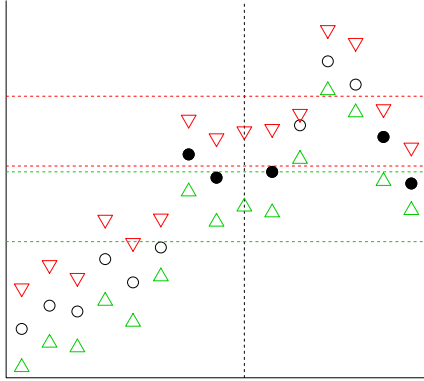
(a) How should we predict the query point's response (dotted line)?

(b) The two primal estimators.



(c) The collective operates.

(d) Predicted value for the query point.



demonstrate that **COBRA** has remarkable speed in terms of CPU timings. In the context of high-dimensional (such as genomic) data, such velocity is critical, and in fact **COBRA** can natively take advantage of multi-core parallel environments.

The paper is organized as follows. In [Section 2](#), we describe the combined estimator—the regression collective—and derive a non-asymptotic risk bound. Next we present the main result, that is the collective is asymptotically at least as good as any of the basic estimators. We also provide a rate of convergence for our procedure which is faster than the usual nonparametric rate.

Section 3 is devoted to the companion R package **COBRA** and presents benchmarks of its excellent performance on both simulated and real data sets, including high-dimensional models. We also show that **COBRA** compares favorably with two competitors, Super Learner (see the seminal paper [van der Laan et al., 2007](#)) and exponentially weighted aggregation (among many other references, see [Dalalyan and Tsybakov, 2008](#)), in that it performs similarly in most situations, much better in some, while it is consistently faster in every case (for the Super Learner). Finally, for ease of exposition, proofs are collected in [Section 4](#).

2 The combined estimator

2.1 Notation

Throughout the article, we assume to be given a training sample denoted by $\mathcal{D}_n = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}$. \mathcal{D}_n is composed of i.i.d. random variables taking their values in $\mathbb{R}^d \times \mathbb{R}$, and distributed as an independent prototype pair (\mathbf{X}, Y) satisfying $\mathbb{E}Y^2 < \infty$ (with the notation $\mathbf{X} = (X_1, \dots, X_d)$). The space \mathbb{R}^d is equipped with the standard Euclidean metric. Our goal is to consistently estimate the regression function $r^*(\mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$, $\mathbf{x} \in \mathbb{R}^d$, using the data \mathcal{D}_n .

To begin with, the original data set \mathcal{D}_n is split into two data sequences $\mathcal{D}_k = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_k, Y_k)\}$ and $\mathcal{D}_\ell = \{(\mathbf{X}_{k+1}, Y_{k+1}), \dots, (\mathbf{X}_n, Y_n)\}$, with $\ell = n - k \geq 1$. For ease of notation, the elements of \mathcal{D}_ℓ are renamed $\{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_\ell, Y_\ell)\}$. There is a slight abuse of notation here, as the same letter is used for both subsets \mathcal{D}_k and \mathcal{D}_ℓ —however, this should not cause any trouble since the context is clear.

Now, suppose that we are given a collection of $M \geq 1$ competing candidates $r_{k,1}, \dots, r_{k,M}$ to estimate r^* . These basic estimators—basic machines—are assumed to be generated using only the first subsample \mathcal{D}_k . These machines can be any among the researcher’s favorite toolkit, such as linear regression, kernel smoother, SVM, Lasso, neural networks, naive Bayes, or random forests. They could equally well be any ad hoc regression rules suggested by the experimental context. The essential idea is that these basic machines can be parametric or nonparametric, or indeed semi-parametric, with possible tuning rules. All what is asked for is that each of the $r_{k,m}(\mathbf{x})$, $m = 1, \dots, M$, is able to provide an estimation of $r^*(\mathbf{x})$ on the basis of \mathcal{D}_k alone. Thus, any collection of model-based or model-free machines are allowed, and the collection is here called the regression collective. Let us emphasize here that the number of basic machines M is considered as fixed throughout the document. Thus the number of machines is not expected to grow and is typically of a reasonable size (M is chosen of the order of 10 in [Section 3](#)).

Given the collection of basic machines $\mathbf{r}_k = (r_{k,1}, \dots, r_{k,M})$, we define the collective estimator T_n to be

$$T_n(\mathbf{r}_k(\mathbf{x})) = \sum_{i=1}^{\ell} W_{n,i}(\mathbf{x}) Y_i, \quad \mathbf{x} \in \mathbb{R}^d,$$

where the random weights $W_{n,i}(\mathbf{x})$ take the form

$$W_{n,i}(\mathbf{x}) = \frac{\mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_\ell\}}}{\sum_{j=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_j)| \leq \varepsilon_\ell\}}}. \quad (2.1)$$

In this definition, ε_ℓ is some positive parameter and, by convention, $0/0 = 0$.

The weighting scheme used in our regression collective is distinctive but not obvious. Starting from [Devroye et al. \(1996\)](#) and [Györfi et al. \(2002\)](#), we see that T_n is a local averaging estimator in the following sense: The value for $r^*(\mathbf{x})$, that is, the estimated outcome at the query point \mathbf{x} , is the unweighted average over those Y_i 's such that \mathbf{X}_i is “close” to the query point. More precisely, for each \mathbf{X}_i in the sample \mathcal{D}_ℓ , “close” means that the output at the query point, generated from each basic machine, is within an ε_ℓ distance of the output generated by the same basic machines at \mathbf{X}_i . If a basic machine evaluated at \mathbf{X}_i is close to the basic machine evaluated at the query point \mathbf{x} , then the corresponding outcome Y_i is included in the average, and not otherwise. Also, as a further note of clarification: “Closeness” of the \mathbf{X}_i is not here to be understood in the Euclidean sense. It refers to closeness of the basic machine outputs at the query point with basic machine outputs over all points in the training data. Training points \mathbf{X}_i 's that are close, in the basic machine sense, to the corresponding basic machine output at the query point contribute to the indicator function for the corresponding outcome Y_i . This discussion is motivated by the fact that a major issue in learning problems consists in building up a metric which is suited to the data (see, *e.g.*, the monograph by [Pekalska and Duin, 2005](#)).

In this context, ε_ℓ plays the role of a smoothing parameter: Put differently, in order to retain Y_i , all basic estimators $r_{k,1}, \dots, r_{k,M}$ have to deliver predictions for the query point \mathbf{x} which are in a ε_ℓ -neighborhood of the predictions $r_{k,1}(\mathbf{X}_i), \dots, r_{k,M}(\mathbf{X}_i)$. Note that the greater ε_ℓ , the more tolerant the process. It turns out that the practical performance of T_n strongly relies on an appropriate choice of ε_ℓ . This important question will thoroughly be discussed in [Section 3](#), where we devise an automatic (*i.e.*, data-dependent) selection strategy of ε_ℓ .

Next, we note that the subscript n in T_n may be a little confusing, since T_n is a weighted average of the Y_i 's in \mathcal{D}_ℓ only. However, T_n depends on the entire data set \mathcal{D}_n , as the rest of the data is used to set up the original machines $r_{k,1}, \dots, r_{k,M}$. Finally, and most importantly, it should be noticed that

the combined estimator T_n is nonlinear with respect to the basic estimators $r_{k,m}$'s. This makes it very different from techniques derived from model selection or convex and linear aggregation literature. As such, it is inspired by the preliminary work of [Mojirsheibani \(1999\)](#) in the supervised classification context.

In addition, let us mention that, in the definition of the weights (2.1), all original estimators are asked to have the same opinion on the importance of the observation \mathbf{X}_i (within the range of ε_ℓ) for the corresponding Y_i to be integrated in the combination T_n . However, this unanimity constraint may be relaxed by imposing, for example, that a fixed fraction $\alpha \in \{1/M, 2/M, \dots, 1\}$ of the machines agree on the importance of \mathbf{X}_i . In that case, the weights take the more sophisticated form

$$W_{n,i}(\mathbf{x}) = \frac{\mathbf{1}_{\{\sum_{m=1}^M \mathbf{1}_{\{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{x}_i)| \leq \varepsilon_\ell}\} \geq M\alpha\}}}{\sum_{j=1}^\ell \mathbf{1}_{\{\sum_{m=1}^M \mathbf{1}_{\{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{x}_j)| \leq \varepsilon_\ell}\} \geq M\alpha\}}}.$$

It turns out that adding the parameter α does not change the asymptotic properties of T_n , provided $\alpha \rightarrow 1$. Thus, to keep a sufficient degree of clarity in the mathematical statements and subsequent proofs, we have decided to consider only the case $\alpha = 1$ (i.e., unanimity). We leave as an exercise the possibility to extend the results to more general values of α . On the other hand, as highlighted by [Section 3](#), α has a nonnegligible impact on the performance of the combined estimator. Accordingly, we will discuss in [Section 3](#) an automatic procedure to select this extra parameter.

2.2 Theoretical performance

This section is devoted to the study of some asymptotic and nonasymptotic properties of the combined estimator T_n , whose quality will be assessed by the quadratic risk

$$\mathbb{E} |T_n(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2.$$

Here and later, \mathbb{E} denotes the expectation with respect to both \mathbf{X} and the sample \mathcal{D}_n . Everywhere in the document, it is assumed that $\mathbb{E}|r_{k,m}(\mathbf{X})|^2 < \infty$ for all $m = 1, \dots, M$. Moreover, we shall need the following technical requirement: For any $m = 1, \dots, M$,

$$r_{k,m}^{-1}((t, +\infty)) \underset{t \uparrow +\infty}{\searrow} \emptyset \quad \text{and} \quad r_{k,m}^{-1}((-\infty, t)) \underset{t \downarrow -\infty}{\searrow} \emptyset. \quad (2.2)$$

It is stressed that this is a mild assumption which is met, for example, whenever the machines are bounded. Throughout, we let

$$T(\mathbf{r}_k(\mathbf{X})) = \mathbb{E}[Y|\mathbf{r}_k(\mathbf{X})]$$

and note that, by the very definition of the L^2 conditional expectation,

$$\mathbb{E} |T(\mathbf{r}_k(\mathbf{X})) - Y|^2 \leq \inf_f \mathbb{E} |f(\mathbf{r}_k(\mathbf{X})) - Y|^2, \quad (2.3)$$

where the infimum is taken over all square integrable functions of $\mathbf{r}_k(\mathbf{X})$.

Our first result is a non-asymptotic inequality, which states that the combined estimator behaves as well as the best one in the original list, within a term measuring how far T_n is from T .

Proposition 2.1. *Let $\mathbf{r}_k = (r_{k,1}, \dots, r_{k,M})$ be the collection of basic estimators, and let $T_n(\mathbf{r}_k(\mathbf{X}))$ be the combined estimator. Then, for all distributions of (\mathbf{X}, Y) with $\mathbb{E}Y^2 < \infty$,*

$$\begin{aligned} \mathbb{E} |T_n(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2 \\ \leq \mathbb{E} |T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 + \inf_f \mathbb{E} |f(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2, \end{aligned}$$

where the infimum is taken over all square integrable functions of $\mathbf{r}_k(\mathbf{X})$. In particular,

$$\begin{aligned} \mathbb{E} |T_n(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2 \\ \leq \min_{m=1, \dots, M} \mathbb{E} |r_{k,m}(\mathbf{X}) - r^*(\mathbf{X})|^2 + \mathbb{E} |T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2. \end{aligned}$$

Note that since, for example, [Proposition 2.1](#) holds for any square integrable function of $\mathbf{r}_k(\mathbf{X})$, this result allows to derive inequalities linking to any existing aggregation procedure: One may consider linear or convex aggregation as well.

[Proposition 2.1](#) reassures us on the performance of T_n with respect to the basic machines, whatever the distribution of (\mathbf{X}, Y) is and regardless of which initial estimator is actually the best. The term $\min_{m=1, \dots, M} \mathbb{E} |r_{k,m}(\mathbf{X}) - r^*(\mathbf{X})|^2$ may be regarded as a bias term, whereas the term $\mathbb{E} |T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2$ is a variance-type term, which can be asymptotically neglected.

Proposition 2.2. *Assume that $\varepsilon_\ell \rightarrow 0$ and $\ell \varepsilon_\ell^M \rightarrow \infty$ as $\ell \rightarrow \infty$. Then*

$$\mathbb{E} |T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 \rightarrow 0 \quad \text{as } \ell \rightarrow \infty,$$

for all distributions of (\mathbf{X}, Y) with $\mathbb{E}Y^2 < \infty$. Thus,

$$\limsup_{\ell \rightarrow \infty} \mathbb{E} |T_n(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2 \leq \min_{m=1, \dots, M} \mathbb{E} |r_{k,m}(\mathbf{X}) - r^*(\mathbf{X})|^2.$$

This result is remarkable, for at least two reasons. Firstly, it shows that, in terms of predictive quadratic risk, the combined estimator does asymptotically at least as well as the best primitive machine. Secondly, the result is universal, in the sense that it is true for all distributions of (\mathbf{X}, Y) , without exceptions.

This is especially interesting because the performance of any estimation procedure eventually depends upon some model and smoothness assumptions on the observations. For example, a linear regression fit performs well if the distribution is truly linear, but may behave poorly otherwise. Similarly, the Lasso procedure is known to do a good job for non-correlated designs (see [van de Geer, 2008](#)), with no clear guarantee however in adversarial situations. Likewise, rates of convergence of nonparametric procedures such as the k -nearest neighbor method, kernel estimators and random forests dramatically deteriorate as the ambient dimension increases, but may be significantly improved if the true underlying dimension is reasonable. This phenomenon is thoroughly analyzed for the random forests algorithm in [Biau \(2012\)](#).

The universal result exhibited in [Proposition 2.2](#) does not require any regularity assumption on the basic machines. However, this universality comes at the price that we have no guarantee on the rate of convergence of the variance term. Nevertheless, assuming some light additional smoothness conditions, one has the following result.

Theorem 2.1. *Assume that Y and the basic machines \mathbf{r}_k are bounded by some constant R . Assume moreover that there exists a constant $L \geq 0$ such that, for every $k \geq 1$,*

$$|T(\mathbf{r}_k(\mathbf{x})) - T(\mathbf{r}_k(\mathbf{y}))| \leq L|\mathbf{r}_k(\mathbf{x}) - \mathbf{r}_k(\mathbf{y})|, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

Then, with the choice $\varepsilon_\ell \propto \ell^{-\frac{1}{M+2}}$, one has

$$\mathbb{E} |T_n(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2 \leq \min_{m=1, \dots, M} \mathbb{E} |r_{k,m}(\mathbf{X}) - r^*(\mathbf{X})|^2 + C\ell^{-\frac{2}{M+2}},$$

for some positive constant $C = C(R, L)$, independent of k .

[Theorem 2.1](#) offers an oracle-type inequality with leading constant 1, stating that the risk of the regression collective is bounded by the lowest risk amongst those of the basic machines, *i.e.*, our procedure mimics the performance of the oracle over the set $\{r_{k,m} : m = 1, \dots, M\}$, plus a remainder term of the order of $\ell^{-2/(M+2)}$ which is the price to pay for aggregating. In our setting, it is important to observe that this term has a limited impact. As a matter of fact, since the number of basic machines M is assumed to be fixed and not too large (the implementation presented in [Section 3](#) considers M at most 6), the remainder term is negligible compared to the standard nonparametric rate $\ell^{-2/(d+2)}$ in dimension d . While the rate $\ell^{-2/(d+2)}$ is affected by the curse of dimensionality when d is large, this is not the case for the term $\ell^{-2/(M+2)}$. Obviously, provided that the distribution of (\mathbf{X}, Y) may be described parametrically and that one of the initial estimators is adapted to this distribution, faster rates of the order of $1/\ell$ could emerge in the bias term. Nonetheless, the regression collective is designed for much more adversarial regression problems, hence the rate exhibited in [Theorem 2.1](#) appears

satisfactory. As a final comment to this result, we stress that our approach carries no assumption on the random design and mild ones over the primal estimators, whereas stringent conditions over the deterministic design and/or the primal estimators are necessary to prove similar results in other aggregation procedures such as the Lasso (Bunea et al., 2007b; van de Geer, 2008).

The crux is that model and smoothness assumptions are usually unverifiable, especially in modern high-dimensional and large scale data sets. To circumvent this difficulty, people often try many different methods and retain the one exhibiting the best empirical results. Our aggregation strategy offers a nice alternative, in the sense that if one of the initial estimators is consistent for a given class \mathcal{M} of distributions, then, under light smoothness assumptions, T_n inherits the same property. To be more precise, assume that the aggregation problem is well-specified, *i.e.*, that one of the original estimators, say r_{k,m_0} , satisfies

$$\mathbb{E} |r_{k,m_0}(\mathbf{X}) - r^*(\mathbf{X})|^2 \rightarrow 0 \quad \text{as } k \rightarrow \infty,$$

for all distribution of (\mathbf{X}, Y) in some class \mathcal{M} . Then, under the assumptions of Theorem 2.1, with the choice $\varepsilon_\ell \propto \ell^{-\frac{1}{M+2}}$, one has

$$\lim_{k, \ell \rightarrow \infty} \mathbb{E} |T_n(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2 = 0.$$

3 Implementation and numerical studies

This section is devoted to the implementation of the described method. Its excellent performance is then assessed in a series of experiments. The companion R package COBRA (standing for COmBined Regression Alternative) is available on the CRAN website <http://cran.r-project.org/web/packages/COBRA/index.html>, for Linux, Mac and Windows platforms, see Guedj (2013). COBRA includes a `parallel` option, allowing for improved performance on multi-core computers (see Knaus, 2010).

As raised in the previous section, a precise calibration of the smoothing parameter ε_ℓ is crucial. Clearly, a too small value will discard many machines and most weights will be zero. Conversely, a large value sets all weights to $1/\Sigma$ with

$$\Sigma = \sum_{j=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{x}_j)| \leq \varepsilon_\ell\}},$$

giving the naive predictor that does not account for any new data point and predicts the mean over the sample \mathcal{D}_ℓ . We also consider a relaxed version of the unanimity constraint: Instead of requiring global agreement over the implemented machines, consider some $\alpha \in (0, 1]$ and keep observation Y_i in

the construction of T_n if and only if at least a proportion α of the machines agree on the importance of \mathbf{X}_i . This parameter requires as well a precise calibration. To understand better, consider the following toy example: On some data set, assume most machines but one have nice predictive performance. For any new data point, requiring global agreement will fail since the pool of machines is heterogeneous. In this regard, α should be seen as a measure of homogeneity: If a small value is selected, it should be seen as an indicator that some machines perform (possibly much) better than some others. Conversely, a large value indicates that the predictive abilities of the machines are close.

A natural measure of the risk in the prediction context is the empirical quadratic loss, namely

$$\hat{R}(\hat{\mathbf{Y}}) = \frac{1}{p} \sum_{j=1}^p (\hat{Y}_j - Y_j)^2,$$

where $\hat{\mathbf{Y}} = (\hat{Y}_1, \dots, \hat{Y}_p)$ is the vector of predicted values for the responses Y_1, \dots, Y_p and $\{(\mathbf{X}_j, Y_j)\}_{j=1}^p$ is a testing sample. We adopted the following protocol: Using a simple data-splitting device, ε_ℓ and α are chosen by minimizing the empirical risk \hat{R} over the set $\{\varepsilon_{\ell, \min}, \dots, \varepsilon_{\ell, \max}\} \times \{1/M, \dots, 1\}$, where $\varepsilon_{\ell, \min} = 10^{-300}$ and $\varepsilon_{\ell, \max}$ is proportional to the largest absolute difference between two predictions of the pool of machines. In the package, $\#\{\varepsilon_{\ell, \min}, \dots, \varepsilon_{\ell, \max}\}$ may be modified by the user, otherwise the default value 200 is chosen. It is also possible to choose either a linear or a logistic scale. [Figure 2](#) illustrates the discussion about the choice of ε_ℓ and α .

By default, COBRA includes the following classical packages dealing with regression estimation and prediction. However, note that the user has the choice to modify this list to her/his own convenience:

- Lasso (R package `lars`, see [Hastie and Efron, 2012](#)).
- Ridge regression (R package `ridge`, see [Cule, 2012](#)).
- k -nearest neighbors (R package `FNN`, see [Li, 2013](#)).
- CART algorithm (R package `tree`, see [Ripley, 2012](#)).
- Random Forests algorithm (R package `randomForest`, see [Liaw and Wiener, 2002](#)).

First, COBRA is benchmarked on synthetic data. For each of the following eight models, two designs are considered: Uniform over $(-1, 1)^d$ (referred to as “Uncorrelated” in [Table 1](#), [Table 2](#) and [Table 3](#)), and Gaussian with mean 0 and covariance matrix Σ with $\Sigma_{ij} = 2^{-|i-j|}$ (“Correlated”). Models considered cover a wide spectrum of contemporary regression problems. Indeed, [Model 1](#)

is a toy example, [Model 2](#) comes from [van der Laan et al. \(2007\)](#), [Model 3](#) and [Model 4](#) appear in [Meier et al. \(2009\)](#). [Model 5](#) is somewhat a classic setting. [Model 6](#) is about predicting labels, [Model 7](#) is inspired by high-dimensional sparse regression problems. Finally, [Model 8](#) deals with probability estimation, linking with nonparametric model-free approaches such as in [Malley et al. \(2012\)](#). In the sequel, we let $\mathcal{N}(\mu, \sigma^2)$ denote a Gaussian random variable with mean μ and variance σ^2 . In the simulations, the training data set was usually set to 80% of the whole sample, then split into two equal parts corresponding to \mathcal{D}_k and \mathcal{D}_ℓ .

Model 1. $n = 800$, $d = 50$, $Y = X_1^2 + \exp(-X_2^2)$.

Model 2. $n = 600$, $d = 100$, $Y = X_1X_2 + X_3^2 - X_4X_7 + X_8X_{10} - X_6^2 + \mathcal{N}(0, 0.5)$.

Model 3. $n = 600$, $d = 100$, $Y = -\sin(2X_1) + X_2^2 + X_3 - \exp(-X_4) + \mathcal{N}(0, 0.5)$.

Model 4. $n = 600$, $d = 100$, $Y = X_1 + (2X_2 - 1)^2 + \sin(2\pi X_3)/(2 - \sin(2\pi X_3)) + \sin(2\pi X_4) + 2\cos(2\pi X_4) + 3\sin^2(2\pi X_4) + 4\cos^2(2\pi X_4) + \mathcal{N}(0, 0.5)$.

Model 5. $n = 700$, $d = 20$, $Y = \mathbf{1}_{\{X_1 > 0\}} + X_2^3 + \mathbf{1}_{\{X_4 + X_6 - X_8 - X_9 > 1 + X_{14}\}} + \exp(-X_2^2) + \mathcal{N}(0, 0.5)$.

Model 6. $n = 500$, $d = 30$, $Y = \sum_{k=1}^{10} \mathbf{1}_{\{X_k^3 < 0\}} - \mathbf{1}_{\{\mathcal{N}(0,1) > 1.25\}}$.

Model 7. $n = 600$, $d = 300$, $Y = X_1^2 + X_2^2X_3 \exp(-|X_4|) + X_6 - X_8 + \mathcal{N}(0, 0.5)$.

Model 8. $n = 600$, $d = 50$, $Y = \mathbf{1}_{\{X_1 + X_4^3 + X_9 + \sin(X_{12}X_{18}) + \mathcal{N}(0,0.1) > 0.38\}}$.

[Table 1](#) presents the mean quadratic error and standard deviation over 100 independent replications, for each model and design. Bold number identifies the lowest error, *i.e.*, the best competitor. Boxplots of errors are presented in [Figure 3](#) and [Figure 4](#). Further, [Figure 5](#) and [Figure 6](#) shows the predictive capacities of COBRA, and [Figure 7](#) depicts its ability to reconstruct the functional dependence over the covariates in the context of additive regression, assessing the striking performance of our approach in a wide spectrum of statistical settings. A remarkable fact is that COBRA performs at least as well as the best machine, and improves even significantly in [Model 3](#), [Model 5](#) and [Model 6](#).

Next, since more and more problems in contemporary statistics involve high-dimensional data, we have tested the abilities of COBRA in that context. As highlighted by [Table 4](#) and [Figure 8](#), the main message is that COBRA is perfectly able to deal with high-dimensional data, provided that it is fed with machines which are known to perform well in such situations (possibly at the price of a sparsity assumption). In that context, we conducted 200 independent replications for the three following models:

Model 9. $n = 500$, $d = 1000$, $Y = X_1 + 3X_3^2 - 2\exp(-X_5) + X_6$. *Uncorrelated design.*

Model 10. $n = 500$, $d = 1000$, $Y = X_1 + 3X_3^2 - 2\exp(-X_5) + X_6$. *Correlated design*

Model 11. $n = 500$, $d = 1500$, $Y = \exp(-X_1) + \exp(X_1) + \sum_{j=2}^d X_j^{j/100}$. *Uncorrelated design.*

A legitimate question that arises is where one should cut the initial sample \mathcal{D}_n ? In other words, for a given data set of size n , what is the optimal value for k ? A naive approach is to cut the initial sample in two halves (i.e., $k = n/2$): This appears to be satisfactory provided that n is large enough, which may be too much of an unrealistic assumption in numerous experimental settings. A more involved choice is to adopt a random cut scheme, where k is chosen uniformly in $\{1, \dots, n\}$. [Figure 9](#) presents the boxplot of errors of the five default machines and COBRA with that random cutting strategy, and also shows the risk of COBRA with respect to k : To illustrate this phenomenon, we tested a thousand random cuts on the following [Model 12](#). As showed in [Figure 9](#), for that particular model, the best value seems to be near $3n/4$.

Model 12. $n = 1200$, $d = 10$, $Y = X_1 + 3X_3^2 - 2\exp(-X_5) + X_6$. *Uncorrelated design.*

The average risk of COBRA on a thousand replications of [Model 12](#) is 0.3124. Since this delivered a thousand prediction vectors, a natural idea is to take their mean or median. The risk of the mean is 0.2306, and the median has an even better risk (0.2184). Since a random cut scheme may generate some unstability, we advise practitioners to compute a few COBRA estimators, then compute the mean or median vector of their predictions.

Next, we compare COBRA to the Super Learner algorithm ([Polley and van der Laan, 2012](#)). This widespread algorithm was first described in [van der Laan et al. \(2007\)](#) and extended in [Polley and van der Laan \(2010\)](#). Super Learner is used in this section as the key competitor to our method. In a nutshell, the Super Learner trains basic machines r_1, \dots, r_M on the whole sample \mathcal{D}_n . Then, following a V -fold cross-validation procedure, Super Learner adopts a V -blocks partition of the set $\{1, \dots, n\}$ and computes the matrix

$$H = (H_{ij})_{1 \leq i \leq n}^{1 \leq j \leq M},$$

where H_{ij} is the prediction for the query point \mathbf{X}_i made by machine j trained on all remaining $V - 1$ blocks, i.e., excluding the block containing \mathbf{X}_i . The Super Learner estimator is then

$$SL = \sum_{j=1}^M \hat{\alpha}_j r_j,$$

where

$$\hat{\alpha} \in \arg \inf_{\alpha \in \Lambda^M} \sum_{i=1}^n |Y_i - (H\alpha)_i|^2,$$

with Λ^M denoting the simplex

$$\Lambda^M = \left\{ \alpha \in \mathbb{R}^M : \sum_{j=1}^M \alpha_j = 1, \alpha_j \geq 0 \text{ for any } j = 1, \dots, M \right\}.$$

Although this convex aggregation scheme is significantly different from our theoretical setting, we feel close to the approach used in the **SuperLearner** package, allowing the user to aggregate as many machines as desired, then blending them to deliver predictive outcomes. For that reason, it is reasonable to deploy Super Learner as a benchmark in our study of regression collectives.

[Table 2](#) summarizes the performance of **COBRA** and **SuperLearner** (used with `SL.randomForest`, `SL.ridge` and `SL.glmnet`, for the fairness of the comparison) through the described protocol. Both methods compete on similar terms in most models, although **COBRA** proves much more efficient on correlated design in [Model 2](#) and [Model 4](#). This already remarkable result is to be stressed by the flexibility and velocity showed by **COBRA**. Indeed, as emphasized in [Table 3](#), without even using the `parallel` option, **COBRA** obtains similar or better results than **SuperLearner** roughly five times faster. Note also that **COBRA** suffers from a disadvantage: **SuperLearner** is built on the whole sample \mathcal{D}_n whereas **COBRA** only uses $\ell < n$ data points. Finally, observe that the algorithmic cost of computing the random weights on n_{test} query points is $\ell \times M \times n_{\text{test}}$ operations. In the package, those calculations are handled in C language for optimal speed performance.

Super Learner is a natural competitor on the implementation side. However, on the theoretical side, it can hardly be considered as a complete benchmark: Thus, we compared **COBRA** to the popular exponentially weighted aggregation method (EWA). We implemented the following version of the EWA: For all preliminary estimators $r_{k,1}, \dots, r_{k,M}$, their empirical risks $\hat{R}_1, \dots, \hat{R}_M$ are computed on a subsample of \mathcal{D}_ℓ and the EWA is

$$\text{EWA}_\beta : \mathbf{x} \mapsto \sum_{j=1}^M \hat{w}_j r_{k,j}(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d,$$

where

$$\hat{w}_j = \frac{\exp(-\beta \hat{R}_j)}{\sum_{i=1}^M \exp(-\beta \hat{R}_i)}, \quad j = 1, \dots, M.$$

The temperature parameter $\beta > 0$ is selected by minimizing the empirical risk of EWA_β over a data-based grid, in the same spirit as the selection of ε_ℓ

and α . We conducted 200 independent replications, on Models 9 to 12. The conclusion is that COBRA outperforms the EWA estimator in some models, and delivers similar performance in others, as shown in Figure 10 and Table 5.

Finally, COBRA is used to process the following real-life data sets:

- Concrete Slump Test⁵ (see Yeh, 2007).
- Concrete Compressive Strength⁶ (see Yeh, 1998).
- Wine Quality⁷ (see Cortez et al., 2009). Note that this data set involves supervised classification and opens a line for future research since COBRA is mainly devoted to regression.

The good predictive performance of COBRA is summarized in Figure 11 and errors are presented in Figure 12. For every data set, the sample is divided into a training set (90%) and a testing set (10%) on which the predictive performance is evaluated. Boxplots are obtained by shuffling the data points a hundred times.

As a conclusion to this thorough experimental protocol, it is our belief that COBRA sets a new gold standard, both in terms of performance and velocity, for prediction-oriented problems in the context of regression.

⁵<http://archive.ics.uci.edu/ml/datasets/Concrete+Slump+Test>.

⁶<http://archive.ics.uci.edu/ml/datasets/Concrete+Compressive+Strength>.

⁷<http://archive.ics.uci.edu/ml/datasets/Wine+Quality>.

Table 1: Quadratic errors of the implemented machines and COBRA. Means and standard deviations over 100 independent replications.

Uncorr.		lars	ridge	fnn	tree	rf	COBRA
Model 1	m.	0.1561	0.1324	0.1585	0.0281	0.0330	0.0259
	sd.	0.0123	0.0094	0.0123	0.0043	0.0033	0.0036
Model 2	m.	0.4880	0.2462	0.3070	0.1746	0.1366	0.1645
	sd.	0.0676	0.0233	0.0303	0.0270	0.0161	0.0207
Model 3	m.	0.2536	0.5347	1.1603	0.4954	0.4027	0.2332
	sd.	0.0271	0.4469	0.1227	0.0772	0.0558	0.0272
Model 4	m.	7.6056	6.3271	10.5890	3.7358	3.5262	3.3640
	sd.	0.9419	1.0800	0.9404	0.8067	0.3223	0.5178
Model 5	m.	0.2943	0.3311	0.5169	0.2918	0.2234	0.2060
	sd.	0.0214	0.1012	0.0439	0.0279	0.0216	0.0210
Model 6	m.	0.8438	1.0303	2.0702	2.3476	1.3354	0.8345
	sd.	0.0916	0.4840	0.2240	0.2814	0.1590	0.1004
Model 7	m.	1.0920	0.5452	0.9459	0.3638	0.3110	0.3052
	sd.	0.2265	0.0920	0.0833	0.0456	0.0325	0.0298
Model 8	m.	0.1308	0.1279	0.2243	0.1715	0.1236	0.1021
	sd.	0.0120	0.0161	0.0189	0.0270	0.0100	0.0155
Corr.		lars	ridge	fnn	tree	rf	COBRA
Model 1	m.	2.3736	1.9785	2.0958	0.3312	0.5766	0.3301
	sd.	0.4108	0.3538	0.3414	0.1285	0.1914	0.1239
Model 2	m.	8.1710	4.0071	4.3892	1.3609	1.4768	1.3612
	sd.	1.5532	0.6840	0.7190	0.4647	0.4415	0.4654
Model 3	m.	6.1448	6.0185	8.2154	4.3175	4.0177	3.7917
	sd.	11.9450	12.0861	13.3121	11.7386	12.4160	11.1806
Model 4	m.	60.5795	42.2117	51.7293	9.6810	14.7731	9.6906
	sd.	11.1303	9.8207	10.9351	3.9807	5.9508	3.9872
Model 5	m.	6.2325	7.1762	10.1254	3.1525	4.2289	2.1743
	sd.	2.4320	3.5448	3.1190	2.1468	2.4826	1.6640
Model 6	m.	1.2765	1.5307	2.5230	2.6185	1.2027	0.9925
	sd.	0.1381	0.9593	0.2762	0.3445	0.1600	0.1210
Model 7	m.	20.8575	4.4367	5.8893	3.6865	2.7318	2.9127
	sd.	7.1821	1.0770	1.2226	1.0139	0.8945	0.9072
Model 8	m.	0.1366	0.1308	0.2267	0.1701	0.1226	0.0984
	sd.	0.0127	0.0143	0.0179	0.0302	0.0102	0.0144

Table 2: Quadratic errors of SuperLearner and COBRA. Means and standard deviations over 100 independent replications.

Uncorr.		SL	COBRA
Model 1	m.	0.0541	0.0320
	sd.	0.0053	0.0104
Model 2	m.	0.1765	0.3569
	sd.	0.0167	0.8797
Model 3	m.	0.2081	0.2573
	sd.	0.0282	0.0699
Model 4	m.	4.3114	3.7464
	sd.	0.4138	0.8746
Model 5	m.	0.2119	0.2187
	sd.	0.0317	0.0427
Model 6	m.	0.7627	1.0220
	sd.	0.1023	0.3347
Model 7	m.	0.1705	0.3103
	sd.	0.0260	0.0490
Model 8	m.	0.1081	0.1075
	sd.	0.0121	0.0235
Corr.		SL	COBRA
Model 1	m.	0.8733	0.3262
	sd.	0.2740	0.1242
Model 2	m.	2.3391	1.3984
	sd.	0.4958	0.3804
Model 3	m.	3.1885	3.3201
	sd.	1.5101	1.8056
Model 4	m.	25.1073	9.3964
	sd.	7.3179	2.8953
Model 5	m.	5.6478	4.9990
	sd.	7.7271	9.3103
Model 6	m.	0.8967	1.1988
	sd.	0.1197	0.4573
Model 7	m.	3.0367	3.1401
	sd.	1.6225	1.6097
Model 8	m.	0.1116	0.1045
	sd.	0.0111	0.0216

Table 3: Average CPU-times in seconds. No parallelization. Means and standard deviations over 10 independent replications.

Uncorr.		SL	COBRA
Model 1	m.	53.92	10.92
	sd.	1.42	0.29
Model 2	m.	57.96	11.90
	sd.	0.95	0.31
Model 3	m.	53.70	10.66
	sd.	0.55	0.11
Model 4	m.	55.00	11.15
	sd.	0.74	0.18
Model 5	m.	28.46	5.01
	sd.	0.73	0.06
Model 6	m.	22.97	3.99
	sd.	0.27	0.05
Model 7	m.	127.80	35.67
	sd.	5.69	1.91
Model 8	m.	32.98	6.46
	sd.	1.33	0.33
Corr.		SL	COBRA
Model 1	m.	61.92	11.96
	sd.	1.85	0.27
Model 2	m.	70.90	14.16
	sd.	2.47	0.57
Model 3	m.	59.91	11.92
	sd.	2.06	0.41
Model 4	m.	63.58	13.11
	sd.	1.21	0.34
Model 5	m.	31.24	5.02
	sd.	0.86	0.07
Model 6	m.	24.29	4.12
	sd.	0.82	0.15
Model 7	m.	145.18	41.28
	sd.	8.97	2.84
Model 8	m.	31.31	6.24
	sd.	0.73	0.11

Table 4: Quadratic errors of the implemented machines and COBRA in high-dimensional situations. Means and standard deviations over 200 independent replications.

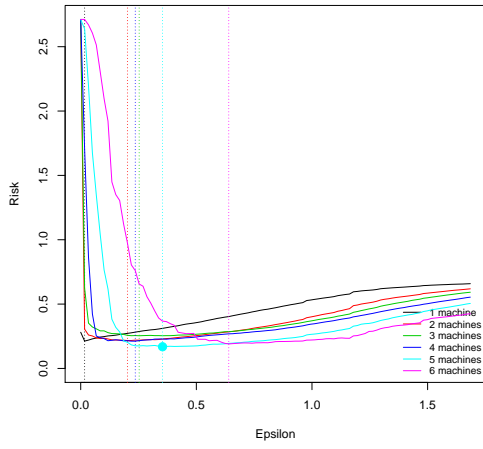
		lars	ridge	fnn	tree	rf	COBRA
Model 9	m.	1.5698	2.9752	3.9285	1.8646	1.5001	0.9996
	sd.	0.2357	0.4171	0.5356	0.3751	0.2491	0.1733
Model 10	m.	5.2356	5.1748	6.1395	6.1585	4.8667	2.7076
	sd.	0.6885	0.7139	0.9192	0.9298	0.6634	0.3810
Model 11	m.	0.1584	0.1055	0.1363	0.0058	0.0327	0.0049
	sd.	0.0199	0.0119	0.0176	0.0010	0.0052	0.0009

Table 5: Quadratic errors of exponentially weighted aggregation (EWA) and COBRA. 200 independent replications.

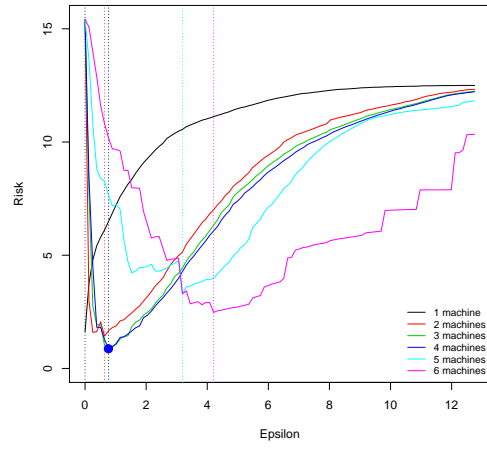
		EWA	COBRA
Model 9	m.	1.1712	1.1360
	sd.	0.2090	0.2468
Model 10	m.	9.4789	12.4353
	sd.	5.6275	9.1267
Model 11	m.	0.0244	0.0128
	sd.	0.0042	0.0237
Model 12	m.	0.4175	0.3124
	sd.	0.0513	0.0884

Figure 2: Examples of calibration of parameters ε_ℓ and α . The bold point is the minimum.

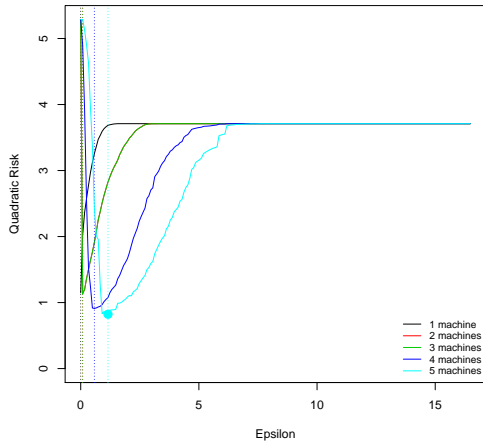
(a) Model 5, uncorrelated design.



(b) Model 5, correlated design.



(c) Model 9.



(d) Model 12.

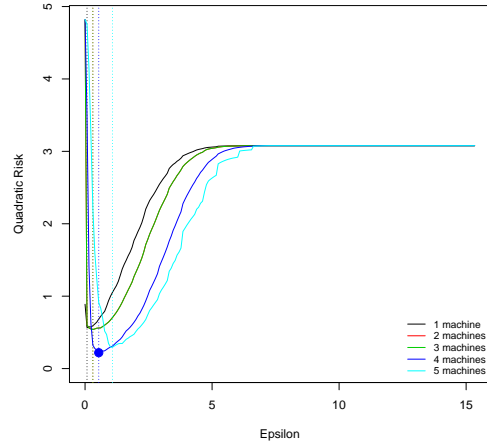


Figure 3: Boxplots of quadratic errors, uncorrelated design. From left to right: lars, ridge, fnn, tree, randomForest, COBRA.

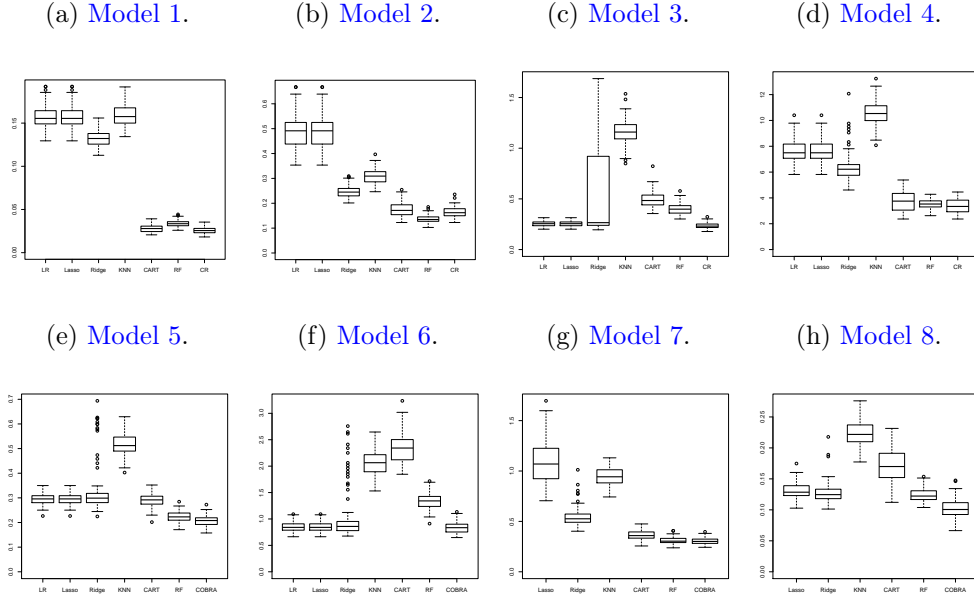


Figure 4: Boxplots of quadratic errors, correlated design. From left to right: lars, ridge, fnn, tree, randomForest, COBRA.

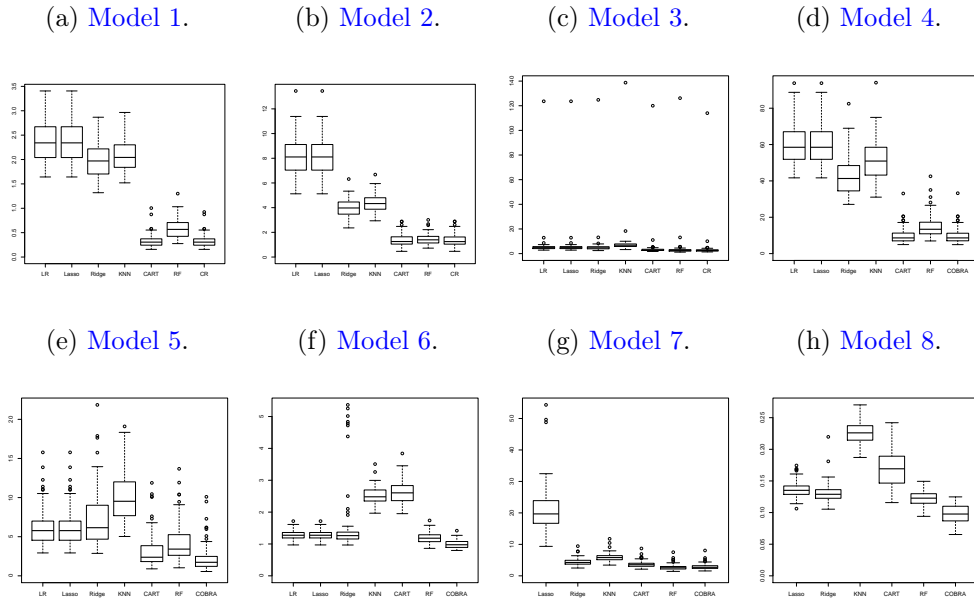


Figure 5: Prediction over the testing set, uncorrelated design. The more points on the first bissectrix, the better the prediction.

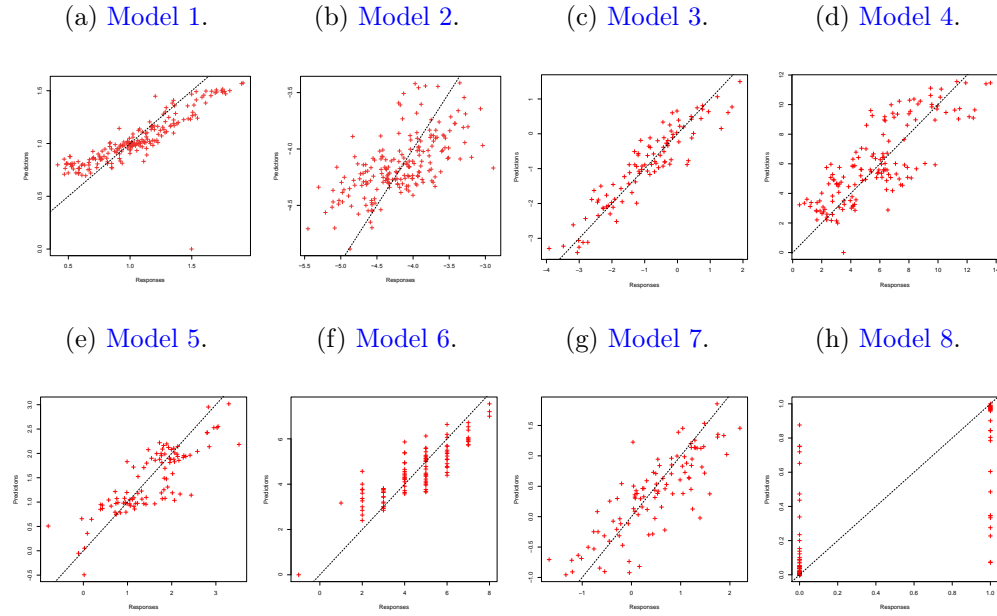


Figure 6: Prediction over the testing set, correlated design. The more points on the first bissectrix, the better the prediction.

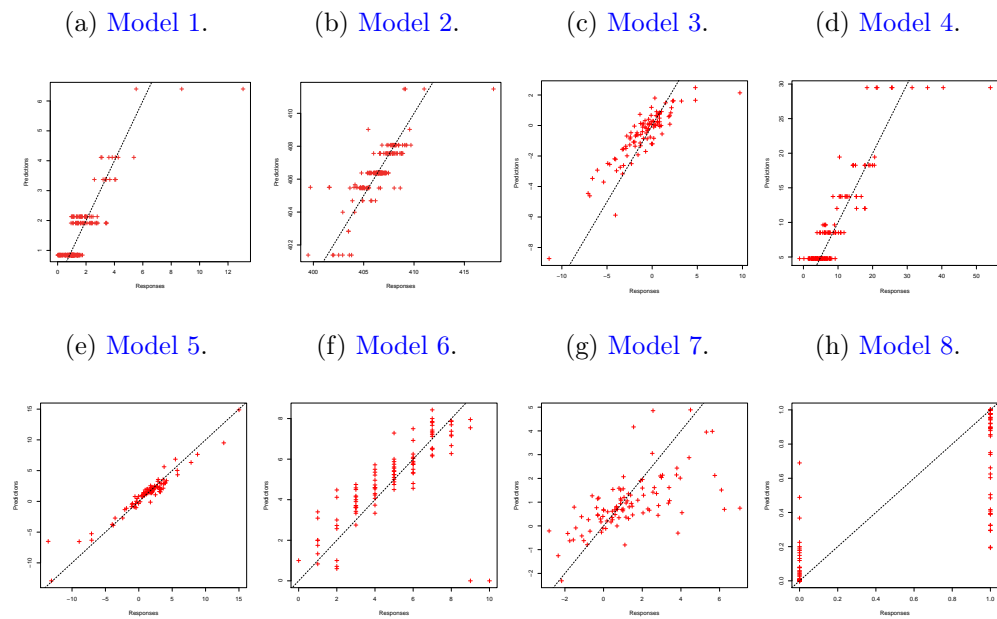
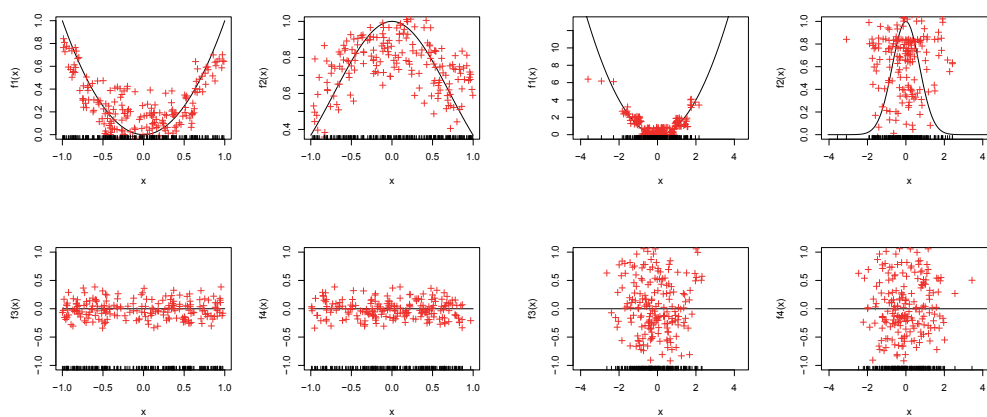


Figure 7: Examples of reconstruction of the functional dependencies, for co-
 variates 1 to 4.

(a) **Model 1**, uncorrelated design.

(b) **Model 1**, correlated design.



(c) **Model 3**, uncorrelated design.

(d) **Model 3**, correlated design.

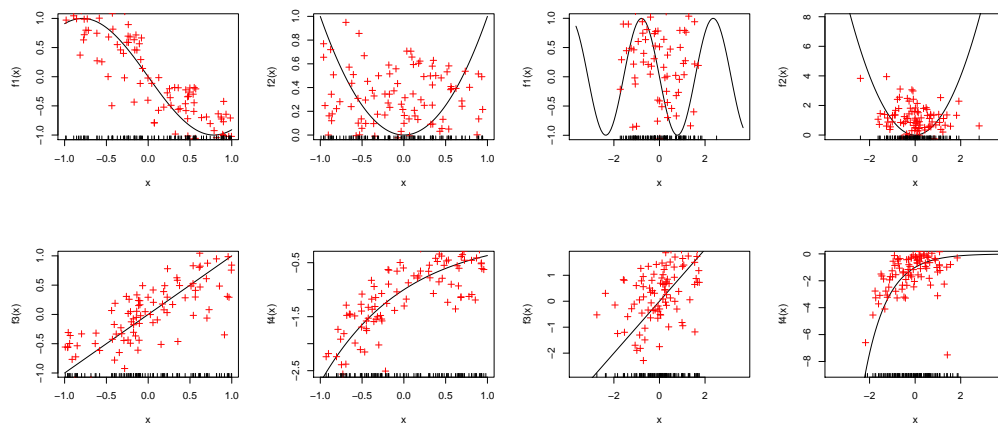


Figure 8: Boxplot of errors, high-dimensional models.

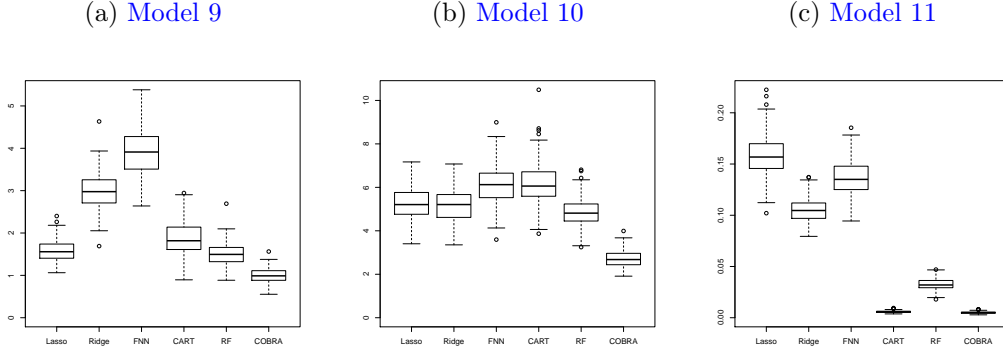


Figure 9: How stable is COBRA?

- (a) Boxplot of errors: Initial sample is randomly cut (1000 replications of Model 12).

(b) Empirical risk with respect to the size of subsample \mathcal{D}_k , in Model 12.

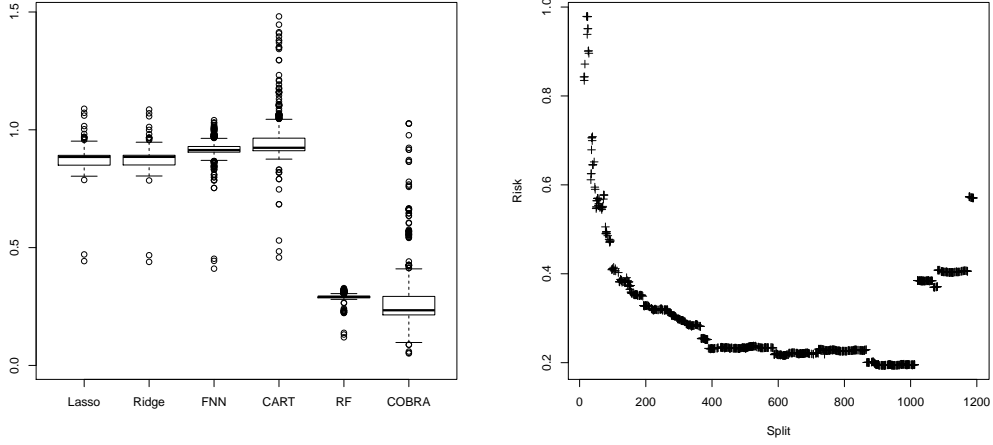


Figure 10: Boxplot of errors: EWA vs COBRA

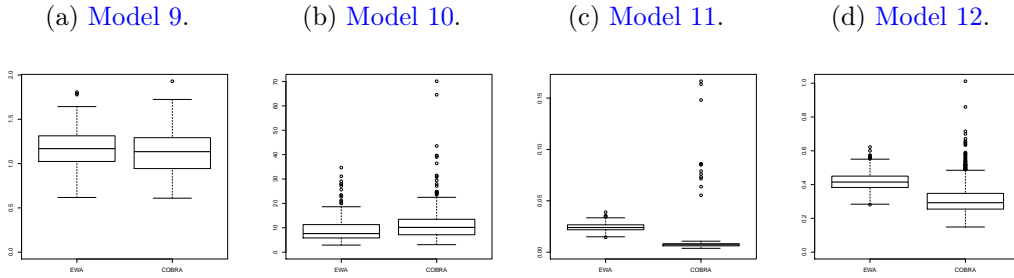
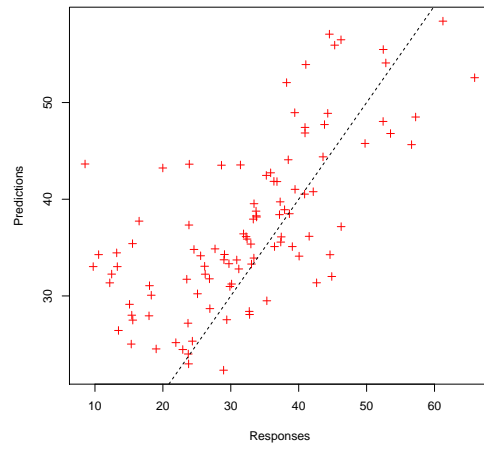
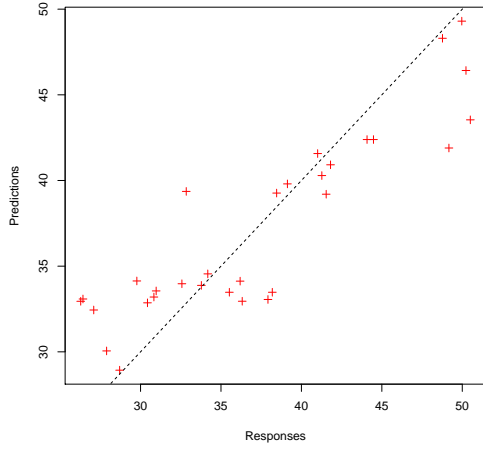


Figure 11: Prediction over the testing set, real-life data sets.

(a) Concrete Slump Test.

(b) Concrete Compressive Strength.



(c) Wine Quality, red wine.

(d) Wine Quality, white wine.

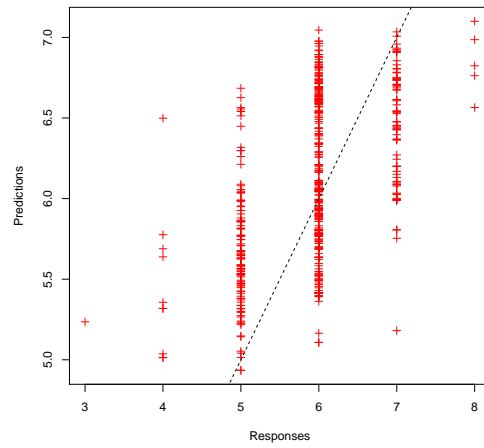
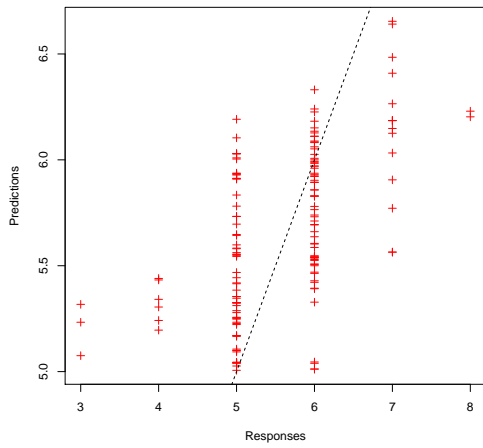


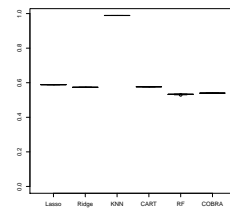
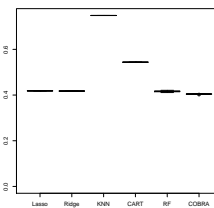
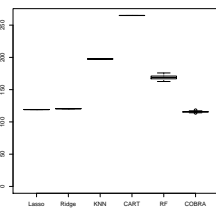
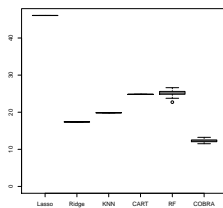
Figure 12: Boxplot of quadratic errors, real-life data sets.

(a) Concrete Slump Test.

(b) Concrete Compressive Strength.

(c) Wine Quality, red wine.

(d) Wine Quality, white wine.



4 Proofs

4.1 Proof of Proposition 2.1

We have

$$\begin{aligned}\mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2 &= \mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 + \mathbb{E}|T(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2 \\ &\quad - 2\mathbb{E}[(T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))(T(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X}))].\end{aligned}$$

As for the double product, notice that

$$\begin{aligned}\mathbb{E}[(T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))(T(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})))] \\ &= \mathbb{E}[\mathbb{E}[(T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))(T(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X}))|\mathbf{r}_k(\mathbf{X}), \mathcal{D}_n]] \\ &= \mathbb{E}[(T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))]\mathbb{E}[T(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|\mathbf{r}_k(\mathbf{X}), \mathcal{D}_n]].\end{aligned}$$

But

$$\begin{aligned}\mathbb{E}[r^*(\mathbf{X})|\mathbf{r}_k(\mathbf{X}), \mathcal{D}_n] &= \mathbb{E}[r^*(\mathbf{X})|\mathbf{r}_k(\mathbf{X})] \\ &\quad (\text{by independence of } \mathbf{X} \text{ and } \mathcal{D}_n) \\ &= \mathbb{E}[\mathbb{E}[Y|\mathbf{X}|\mathbf{r}_k(\mathbf{X})] \\ &= \mathbb{E}[Y|\mathbf{r}_k(\mathbf{X})] \\ &\quad (\text{since } \sigma(\mathbf{r}_k(\mathbf{X})) \subset \sigma(\mathbf{X})) \\ &= T(\mathbf{r}_k(\mathbf{X})).\end{aligned}$$

Consequently,

$$\mathbb{E}[(T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))(T(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})))] = 0$$

and

$$\mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2 = \mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 + \mathbb{E}|T(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2.$$

Thus, by definition of the conditional expectation, and using the fact that $T(\mathbf{r}_k(\mathbf{X})) = \mathbb{E}[r^*(\mathbf{X})|\mathbf{r}_k(\mathbf{X})]$,

$$\mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2 \leq \mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 + \inf_f \mathbb{E}|f(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2,$$

where the infimum is taken over all square integrable functions of $\mathbf{r}_k(\mathbf{X})$. In particular,

$$\begin{aligned}\mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - r^*(\mathbf{X})|^2 \\ \leq \min_{m=1, \dots, M} \mathbb{E}|r_{k,m}(\mathbf{X}) - r^*(\mathbf{X})|^2 + \mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2,\end{aligned}$$

as desired.

4.2 Proof of Proposition 2.2

Note that the second statement is an immediate consequence of the first statement and Proposition 2.1, therefore we only have to prove that

$$\mathbb{E} |T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 \rightarrow 0 \quad \text{as } \ell \rightarrow \infty.$$

We start with a technical lemma, whose proof can be found in the monograph by Györfi et al. (2002).

Lemma 4.1. *Let $B(n, p)$ be a binomial random variable with parameters $n \geq 1$ and $p > 0$. Then*

$$\mathbb{E} \left[\frac{1}{1 + B(n, p)} \right] \leq \frac{1}{p(n+1)}$$

and

$$\mathbb{E} \left[\frac{\mathbf{1}_{\{B(n, p) > 0\}}}{B(n, p)} \right] \leq \frac{2}{p(n+1)}.$$

For all distribution of (\mathbf{X}, Y) , using the elementary inequality $(a + b + c)^2 \leq 3(a^2 + b^2 + c^2)$, note that

$$\begin{aligned} & \mathbb{E} |T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 \\ &= \mathbb{E} \left| \sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) (Y_i - T(\mathbf{r}_k(\mathbf{X}_i)) + T(\mathbf{r}_k(\mathbf{X}_i)) - T(\mathbf{r}_k(\mathbf{X})) + T(\mathbf{r}_k(\mathbf{X}))) \right. \\ & \quad \left. - T(\mathbf{r}_k(\mathbf{X})) \right|^2 \\ &\leq 3\mathbb{E} \left| \sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) (T(\mathbf{r}_k(\mathbf{X}_i)) - T(\mathbf{r}_k(\mathbf{X}))) \right|^2 \end{aligned} \tag{4.1}$$

$$+ 3\mathbb{E} \left| \sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) (Y_i - T(\mathbf{r}_k(\mathbf{X}_i))) \right|^2 \tag{4.2}$$

$$+ 3\mathbb{E} \left| \left(\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) - 1 \right) T(\mathbf{r}_k(\mathbf{X})) \right|^2. \tag{4.3}$$

Consequently, to prove the proposition, it suffices to establish that (4.1), (4.2) and (4.3) tend to 0 as ℓ tends to infinity. This is done, respectively, in Proposition 4.1, Proposition 4.2 and Proposition 4.3 below.

Proposition 4.1. *Under the assumptions of Proposition 2.2,*

$$\lim_{\ell \rightarrow \infty} \mathbb{E} \left| \sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) (T(\mathbf{r}_k(\mathbf{X}_i)) - T(\mathbf{r}_k(\mathbf{X}))) \right|^2 = 0.$$

Proof of Proposition 4.1. By the Cauchy-Schwarz inequality,

$$\begin{aligned}
& \mathbb{E} \left| \sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) (T(\mathbf{r}_k(\mathbf{X}_i)) - T(\mathbf{r}_k(\mathbf{X}))) \right|^2 \\
&= \mathbb{E} \left| \sum_{i=1}^{\ell} \sqrt{W_{n,i}(\mathbf{X})} \sqrt{W_{n,i}(\mathbf{X})} (T(\mathbf{r}_k(\mathbf{X}_i)) - T(\mathbf{r}_k(\mathbf{X}))) \right|^2 \\
&\leq \mathbb{E} \left[\sum_{j=1}^{\ell} W_{n,j}(\mathbf{X}) \sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |T(\mathbf{r}_k(\mathbf{X}_i)) - T(\mathbf{r}_k(\mathbf{X}))|^2 \right] \\
&= \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |T(\mathbf{r}_k(\mathbf{X}_i)) - T(\mathbf{r}_k(\mathbf{X}))|^2 \right] \\
&:= A_n.
\end{aligned}$$

The function T is such that $\mathbb{E}[T^2(\mathbf{r}_k(\mathbf{X}))] < \infty$. Therefore, it can be approximated in an L^2 sense by a continuous function with compact support, say \tilde{T} . This result may be found in many references, amongst them Györfi et al. (2002, Theorem A.1). More precisely, for any $\eta > 0$, there exists a function \tilde{T} such that

$$\mathbb{E} \left| T(\mathbf{r}_k(\mathbf{X})) - \tilde{T}(\mathbf{r}_k(\mathbf{X})) \right|^2 < \eta.$$

Consequently, we obtain

$$\begin{aligned}
A_n &= \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |T(\mathbf{r}_k(\mathbf{X}_i)) - T(\mathbf{r}_k(\mathbf{X}))|^2 \right] \\
&\leq 3\mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |T(\mathbf{r}_k(\mathbf{X}_i)) - \tilde{T}(\mathbf{r}_k(\mathbf{X}_i))|^2 \right] \\
&\quad + 3\mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |\tilde{T}(\mathbf{r}_k(\mathbf{X}_i)) - \tilde{T}(\mathbf{r}_k(\mathbf{X}))|^2 \right] \\
&\quad + 3\mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |\tilde{T}(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 \right] \\
&:= 3A_{n1} + 3A_{n2} + 3A_{n3}.
\end{aligned}$$

Computation of A_{n3} . Thanks to the approximation of T by \tilde{T} ,

$$\begin{aligned}
A_{n3} &= \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |T(\mathbf{r}_k(\mathbf{X})) - \tilde{T}(\mathbf{r}_k(\mathbf{X}))|^2 \right] \\
&\leq \mathbb{E} \left| T(\mathbf{r}_k(\mathbf{X})) - \tilde{T}(\mathbf{r}_k(\mathbf{X})) \right|^2 < \eta.
\end{aligned}$$

Computation of A_{n1} . Denote by μ the distribution of \mathbf{X} . Then,

$$\begin{aligned}
A_{n1} &= \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |\tilde{T}(\mathbf{r}_k(\mathbf{X}_i)) - T(\mathbf{r}_k(\mathbf{X}_i))|^2 \right] \\
&= \ell \mathbb{E} \left[\frac{\mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_1)| \leq \varepsilon_\ell\}}}{\sum_{j=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_j)| \leq \varepsilon_\ell\}}} |\tilde{T}(\mathbf{r}_k(\mathbf{X}_1)) - T(\mathbf{r}_k(\mathbf{X}_1))|^2 \right] \\
&= \ell \mathbb{E} \left\{ \int_{\mathbb{R}^d} |\tilde{T}(\mathbf{r}_k(\mathbf{u})) - T(\mathbf{r}_k(\mathbf{u}))|^2 \right. \\
&\quad \times \mathbb{E} \left[\int_{\mathbb{R}^d} \frac{\mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{u})| \leq \varepsilon_\ell\}}}{\mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{u})| \leq \varepsilon_\ell\}} + \sum_{j=2}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_j)| \leq \varepsilon_\ell\}}} \mu(d\mathbf{x}) \middle| \mathcal{D}_k \right] \\
&\quad \left. \mu(d\mathbf{u}) \right\}.
\end{aligned}$$

Let us prove that

$$\begin{aligned}
A'_{n1} &= \mathbb{E} \left[\int_{\mathbb{R}^d} \frac{\mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{u})| \leq \varepsilon_\ell\}}}{\mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{u})| \leq \varepsilon_\ell\}} + \sum_{j=2}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_j)| \leq \varepsilon_\ell\}}} \mu(d\mathbf{x}) \right. \\
&\quad \left. \middle| \mathcal{D}_k \right] \\
&\leq \frac{2^M}{\ell}.
\end{aligned}$$

To this aim, observe that

$$\begin{aligned}
A'_{n1} &= \mathbb{E} \left[\int_{\mathbb{R}^d} \frac{\mathbf{1}_{\{\mathbf{x} \in \cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{u}) - \varepsilon_\ell, r_{k,m}(\mathbf{u}) + \varepsilon_\ell])\}}}{1 + \sum_{j=2}^{\ell} \mathbf{1}_{\{\mathbf{x}_j \in \cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{x}) - \varepsilon_\ell, r_{k,m}(\mathbf{x}) + \varepsilon_\ell])\}}} \mu(d\mathbf{x}) \middle| \mathcal{D}_k \right] \\
&= \mathbb{E} \left[\int_{\mathbb{R}^d} \frac{\mathbf{1}_{\{\mathbf{x} \in \cup_{(a_1, \dots, a_M) \in \{1,2\}^M} r_{k,1}^{-1}(I_{n,1}^{a_1}(\mathbf{u})) \cap \dots \cap r_{k,M}^{-1}(I_{n,M}^{a_M}(\mathbf{u}))\}}}{1 + \sum_{j=2}^{\ell} \mathbf{1}_{\{\mathbf{x}_j \in \cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{x}) - \varepsilon_\ell, r_{k,m}(\mathbf{x}) + \varepsilon_\ell])\}}} \mu(d\mathbf{x}) \middle| \mathcal{D}_k \right] \\
&\leq \sum_{p=1}^{2^M} \mathbb{E} \left[\int_{\mathbb{R}^d} \frac{\mathbf{1}_{\{\mathbf{x} \in R_n^p(\mathbf{u})\}}}{1 + \sum_{j=2}^{\ell} \mathbf{1}_{\{\mathbf{x}_j \in \cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{x}) - \varepsilon_\ell, r_{k,m}(\mathbf{x}) + \varepsilon_\ell])\}}} \mu(d\mathbf{x}) \middle| \mathcal{D}_k \right].
\end{aligned}$$

Here, $I_{n,m}^1(\mathbf{u}) = [r_{k,m}(\mathbf{u}) - \varepsilon_\ell, r_{k,m}(\mathbf{u})]$, $I_{n,m}^2(\mathbf{u}) = [r_{k,m}(\mathbf{u}), r_{k,m}(\mathbf{u}) + \varepsilon_\ell]$, and $R_n^p(\mathbf{u})$ is the p -th set of the form $r_{k,1}^{-1}(I_{n,1}^{a_1}(\mathbf{u})) \cap \dots \cap r_{k,M}^{-1}(I_{n,M}^{a_M}(\mathbf{u}))$ assuming that they have been ordered using the lexicographic order of (a_1, \dots, a_M) .

Next, note that

$$\mathbf{x} \in R_n^p(\mathbf{u}) \Rightarrow R_n^p(\mathbf{u}) \subset \bigcap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{x}) - \varepsilon_\ell, r_{k,m}(\mathbf{x}) + \varepsilon_\ell]).$$

To see this, just observe that, for all $m = 1, \dots, M$, if $r_{k,m}(\mathbf{z}) \in [r_{k,m}(\mathbf{u}) - \varepsilon_\ell, r_{k,m}(\mathbf{u})]$, i.e., $r_{k,m}(\mathbf{u}) - \varepsilon_\ell \leq r_{k,m}(\mathbf{z}) \leq r_{k,m}(\mathbf{u})$, then, as $r_{k,m}(\mathbf{u}) - \varepsilon_\ell \leq r_{k,m}(\mathbf{x}) \leq r_{k,m}(\mathbf{u})$, one has $r_{k,m}(\mathbf{x}) - \varepsilon_\ell \leq r_{k,m}(\mathbf{z}) \leq r_{k,m}(\mathbf{x}) + \varepsilon_\ell$. Similarly, if $r_{k,m}(\mathbf{u}) \leq r_{k,m}(\mathbf{z}) \leq r_{k,m}(\mathbf{u}) + \varepsilon_\ell$, then $r_{k,m}(\mathbf{u}) \leq r_{k,m}(\mathbf{x}) \leq r_{k,m}(\mathbf{u}) + \varepsilon_\ell$ implies $r_{k,m}(\mathbf{x}) - \varepsilon_\ell \leq r_{k,m}(\mathbf{z}) \leq r_{k,m}(\mathbf{x}) + \varepsilon_\ell$. Consequently,

$$\begin{aligned} A'_{n1} &\leq \sum_{p=1}^{2^M} \mathbb{E} \left[\int_{\mathbb{R}^d} \frac{\mathbf{1}_{\{\mathbf{x} \in R_n^p(\mathbf{u})\}}}{1 + \sum_{j=2}^{\ell} \mathbf{1}_{\{\mathbf{x}_j \in R_n^p(\mathbf{u})\}}} \mu(d\mathbf{x}) \middle| \mathcal{D}_k \right] \\ &= \sum_{p=1}^{2^M} \mathbb{E} \left[\frac{\mu\{R_n^p(\mathbf{u})\}}{1 + \sum_{j=2}^{\ell} \mathbf{1}_{\{\mathbf{x}_j \in R_n^p(\mathbf{u})\}}} \middle| \mathcal{D}_k \right] \\ &\leq \sum_{p=1}^{2^M} \mathbb{E} \left[\frac{\mu\{R_n^p(\mathbf{u})\}}{\ell \mu\{R_n^p(\mathbf{u})\}} \middle| \mathcal{D}_k \right] \\ &\leq \frac{2^M}{\ell} \end{aligned}$$

(by the first statement of [Lemma 4.1](#)). Thus, returning to A_{n1} , we obtain

$$A_{n1} \leq 2^M \mathbb{E} \left| \tilde{T}(\mathbf{r}_k(\mathbf{X}) - T(\mathbf{r}_k(\mathbf{X}))) \right|^2 < 2^M \eta.$$

Computation of A_{n2} . For any $\delta > 0$, write

$$\begin{aligned} A_{n2} &= \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |\tilde{T}(\mathbf{r}_k(\mathbf{X}_i)) - \tilde{T}(\mathbf{r}_k(\mathbf{X}))|^2 \right] \\ &= \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |\tilde{T}(\mathbf{r}_k(\mathbf{X}_i)) - \tilde{T}(\mathbf{r}_k(\mathbf{X}))|^2 \mathbf{1}_{\bigcup_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_i)| > \delta\}} \right] \\ &\quad + \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |\tilde{T}(\mathbf{r}_k(\mathbf{X}_i)) - \tilde{T}(\mathbf{r}_k(\mathbf{X}))|^2 \mathbf{1}_{\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_i)| \leq \delta\}} \right] \\ &\leq 4 \sup_{\mathbf{u} \in \mathbb{R}^d} |\tilde{T}(\mathbf{r}_k(\mathbf{u}))|^2 \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) \mathbf{1}_{\bigcup_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_i)| > \delta\}} \right] \quad (4.4) \end{aligned}$$

$$+ \left(\sup_{\mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \bigcap_{m=1}^M \{|r_{k,m}(\mathbf{u}) - r_{k,m}(\mathbf{v})| \leq \delta\}} |\tilde{T}(\mathbf{r}_k(\mathbf{v})) - \tilde{T}(\mathbf{r}_k(\mathbf{u}))| \right)^2. \quad (4.5)$$

With respect to the term (4.4), if $\delta > \varepsilon_\ell$, then

$$\begin{aligned} &\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) \mathbf{1}_{\bigcup_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_i)| > \delta\}} \\ &= \sum_{i=1}^{\ell} \frac{\mathbf{1}_{\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_\ell\}} \mathbf{1}_{\bigcup_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_i)| > \delta\}}}{\sum_{j=1}^{\ell} \mathbf{1}_{\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_j)| \leq \varepsilon_\ell\}}} \\ &= 0. \end{aligned}$$

It follows that, for all $\delta > 0$, this term converges to 0 as ℓ tends to infinity. On the other hand, letting $\delta \rightarrow 0$, we see that the term (4.5) tends to 0 as well, by uniform continuity of \tilde{T} . Hence, A_{n2} tends to 0 as ℓ tends to infinity. Letting finally η go to 0, we conclude that A_n vanishes as ℓ tends to infinity. \square

Proposition 4.2. *Under the assumptions of Proposition 2.2,*

$$\lim_{\ell \rightarrow \infty} \mathbb{E} \left| \sum_{i=1}^{\ell} W_{n,i}(\mathbf{X})(Y_i - T(\mathbf{r}_k(\mathbf{X}_i))) \right|^2 = 0.$$

Proof of Proposition 4.2.

$$\begin{aligned} & \mathbb{E} \left| \sum_{i=1}^{\ell} W_{n,i}(\mathbf{X})(Y_i - T(\mathbf{r}_k(\mathbf{X}_i))) \right|^2 \\ &= \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \mathbb{E}[W_{n,i}(\mathbf{X})W_{n,j}(\mathbf{X})(Y_i - T(\mathbf{r}_k(\mathbf{X}_i)))(Y_j - T(\mathbf{r}_k(\mathbf{X}_j)))] \\ &= \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}^2(\mathbf{X}) |Y_i - T(\mathbf{r}_k(\mathbf{X}_i))|^2 \right] \\ &= \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}^2(\mathbf{X}) \sigma^2(\mathbf{r}_k(\mathbf{X}_i)) \right], \end{aligned}$$

where

$$\sigma^2(\mathbf{r}_k(\mathbf{x})) = \mathbb{E}[|Y - T(\mathbf{r}_k(\mathbf{X}))|^2 | \mathbf{r}_k(\mathbf{x})].$$

For any $\eta > 0$, using again Györfi et al. (2002, Theorem A.1), σ^2 can be approximated in an L^1 sense by a continuous function with compact support $\tilde{\sigma}^2$, i.e.,

$$\mathbb{E}|\tilde{\sigma}^2(\mathbf{r}_k(\mathbf{X})) - \sigma^2(\mathbf{r}_k(\mathbf{X}))| < \eta.$$

Thus

$$\begin{aligned} & \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}^2(\mathbf{X}) \sigma^2(\mathbf{r}_k(\mathbf{X}_i)) \right] \\ & \leq \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}^2(\mathbf{X}) \tilde{\sigma}^2(\mathbf{r}_k(\mathbf{X}_i)) \right] \\ & \quad + \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}^2(\mathbf{X}) |\sigma^2(\mathbf{r}_k(\mathbf{X}_i)) - \tilde{\sigma}^2(\mathbf{r}_k(\mathbf{X}_i))| \right] \\ & \leq \sup_{\mathbf{u} \in \mathbb{R}^d} |\tilde{\sigma}^2(\mathbf{r}_k(\mathbf{u}))| \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}^2(\mathbf{X}) \right] \\ & \quad + \mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |\sigma^2(\mathbf{r}_k(\mathbf{X}_i)) - \tilde{\sigma}^2(\mathbf{r}_k(\mathbf{X}_i))| \right]. \end{aligned}$$

With the same argument as for A_{n1} , we obtain

$$\mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) |\sigma^2(\mathbf{r}_k(\mathbf{X}_i)) - \tilde{\sigma}^2(\mathbf{r}_k(\mathbf{X}_i))| \right] \leq 2^M \eta.$$

Therefore, it remains to prove that $\mathbb{E} \left[\sum_{i=1}^{\ell} W_{n,i}^2(\mathbf{X}) \right] \rightarrow 0$ as $\ell \rightarrow \infty$. To this aim, fix $\delta > 0$, and note that

$$\begin{aligned} \sum_{i=1}^{\ell} W_{n,i}^2(\mathbf{X}) &= \frac{\sum_{i=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_{\ell}\}}}{\left(\sum_{j=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_j)| \leq \varepsilon_{\ell}\}} \right)^2} \\ &\leq \min \left\{ \delta, \frac{1}{\sum_{i=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_{\ell}\}}} \right\} \\ &\leq \delta + \frac{\mathbf{1}_{\left\{ \sum_{i=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_{\ell}\}} > 0 \right\}}}{\sum_{i=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_{\ell}\}}}. \end{aligned}$$

To complete the proof, we have to establish that the expectation of the right-hand term tends to 0. Denoting by I a bounded interval on the real line, we have

$$\begin{aligned} &\mathbb{E} \left[\frac{\mathbf{1}_{\left\{ \sum_{i=1}^{\ell} \mathbf{1}_{\{\mathbf{X}_i \in \cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{X}) - \varepsilon_{\ell}, r_{k,m}(\mathbf{X}) + \varepsilon_{\ell}])\}} > 0 \right\}}}{\sum_{i=1}^{\ell} \mathbf{1}_{\{\mathbf{X}_i \in \cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{X}) - \varepsilon_{\ell}, r_{k,m}(\mathbf{X}) + \varepsilon_{\ell}])\}}} \right] \\ &\leq \mathbb{E} \left[\frac{\mathbf{1}_{\left\{ \sum_{i=1}^{\ell} \mathbf{1}_{\{\mathbf{X}_i \in \cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{X}) - \varepsilon_{\ell}, r_{k,m}(\mathbf{X}) + \varepsilon_{\ell}])\}} > 0 \right\}} \mathbf{1}_{\{\mathbf{X} \in \cap_{m=1}^M r_{k,m}^{-1}(I)\}}}{\sum_{i=1}^{\ell} \mathbf{1}_{\{\mathbf{X}_i \in \cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{X}) - \varepsilon_{\ell}, r_{k,m}(\mathbf{X}) + \varepsilon_{\ell}])\}}} \right] \\ &\quad + \mu \left(\bigcup_{m=1}^M r_{k,m}^{-1}(I^c) \right) \\ &= \mathbb{E} \left[\mathbb{E} \left[\frac{\mathbf{1}_{\left\{ \sum_{i=1}^{\ell} \mathbf{1}_{\{\mathbf{X}_i \in \cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{X}) - \varepsilon_{\ell}, r_{k,m}(\mathbf{X}) + \varepsilon_{\ell}])\}} > 0 \right\}} \mathbf{1}_{\{\mathbf{X} \in \cap_{m=1}^M r_{k,m}^{-1}(I)\}}}{\sum_{i=1}^{\ell} \mathbf{1}_{\{\mathbf{X}_i \in \cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{X}) - \varepsilon_{\ell}, r_{k,m}(\mathbf{X}) + \varepsilon_{\ell}])\}}} \right. \right. \\ &\quad \left. \left. \middle| \mathcal{D}_k, \mathbf{X} \right] \right] + \mu \left(\bigcup_{m=1}^M r_{k,m}^{-1}(I^c) \right) \\ &\leq \frac{2}{(\ell + 1)} \mathbb{E} \left[\frac{\mathbf{1}_{\{\mathbf{X} \in \cap_{m=1}^M r_{k,m}^{-1}(I)\}}}{\mu(\cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{X}) - \varepsilon_{\ell}, r_{k,m}(\mathbf{X}) + \varepsilon_{\ell}]))} \right] \\ &\quad + \mu \left(\bigcup_{m=1}^M r_{k,m}^{-1}(I^c) \right). \end{aligned}$$

The last inequality arises from the second statement of [Lemma 4.1](#). By an appropriate choice of I , according to the technical statement [\(2.2\)](#), the second term on the right-hand side can be made as small as desired. Regarding the first term, there exists a finite number N_ℓ of points $\mathbf{z}_1, \dots, \mathbf{z}_{N_\ell}$ such that

$$\bigcap_{m=1}^M r_{k,m}^{-1}(I) \subset \bigcup_{(j_1, \dots, j_M) \in \{1, \dots, N_\ell\}^M} r_{k,1}^{-1}(I_{n,1}(\mathbf{z}_{j_1})) \cap \dots \cap r_{k,M}^{-1}(I_{n,M}(\mathbf{z}_{j_M})),$$

where $I_{n,m}(\mathbf{z}_j) = [\mathbf{z}_j - \varepsilon_\ell/2, \mathbf{z}_j + \varepsilon_\ell/2]$. Suppose, without loss of generality, that the sets

$$r_{k,1}^{-1}(I_{n,1}(\mathbf{z}_{j_1})) \cap \dots \cap r_{k,M}^{-1}(I_{n,M}(\mathbf{z}_{j_M}))$$

are ordered, and denote by R_n^p the p -th among the $N_\ell^M = (\lceil |I|/\varepsilon_\ell \rceil)^M$ sets. Here $|I|$ denotes the length of the interval I and $\lceil x \rceil$ denotes the smallest integer greater than x . For all p ,

$$\mathbf{x} \in R_n^p \Rightarrow R_n^p \subset \bigcap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{x}) - \varepsilon_\ell, r_{k,m}(\mathbf{x}) + \varepsilon_\ell]).$$

Indeed, if $\mathbf{v} \in R_n^p$, then, for all $m = 1, \dots, M$, there exists $j \in \{1, \dots, N_\ell\}$ such that $r_{k,m}(\mathbf{v}) \in [\mathbf{z}_j - \varepsilon_\ell/2, \mathbf{z}_j + \varepsilon_\ell/2]$, that is $\mathbf{z}_j - \varepsilon_\ell/2 \leq r_{k,m}(\mathbf{v}) \leq \mathbf{z}_j + \varepsilon_\ell/2$. Since we also have $\mathbf{z}_j - \varepsilon_\ell/2 \leq r_{k,m}(\mathbf{X}) \leq \mathbf{z}_j + \varepsilon_\ell/2$, we obtain $r_{k,m}(\mathbf{X}) - \varepsilon_\ell \leq r_{k,m}(\mathbf{v}) \leq r_{k,m}(\mathbf{X}) + \varepsilon_\ell$. In conclusion,

$$\begin{aligned} & \mathbb{E} \left[\frac{\mathbf{1}_{\{\mathbf{x} \in \bigcap_{m=1}^M r_{k,m}^{-1}(I)\}}}{\mu(\bigcap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{X}) - \varepsilon_\ell, r_{k,m}(\mathbf{X}) + \varepsilon_\ell]))} \right] \\ & \leq \sum_{p=1}^{N_\ell^M} \mathbb{E} \left[\frac{\mathbf{1}_{\{\mathbf{x} \in R_n^p\}}}{\mu(\bigcap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{X}) - \varepsilon_\ell, r_{k,m}(\mathbf{X}) + \varepsilon_\ell]))} \right] \\ & \leq \sum_{p=1}^{N_\ell^M} \mathbb{E} \left[\frac{\mathbf{1}_{\{\mathbf{x} \in R_n^p\}}}{\mu(R_n^p)} \right] \\ & = N_\ell^M \\ & = \left\lceil \frac{|I|}{\varepsilon_\ell} \right\rceil^M. \end{aligned}$$

The result follows from the assumption $\lim_{\ell \rightarrow \infty} \ell \varepsilon_\ell^M = \infty$. □

Proposition 4.3. *Under the assumptions of [Proposition 2.2](#),*

$$\lim_{\ell \rightarrow \infty} \mathbb{E} \left| \left(\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) - 1 \right) T(\mathbf{r}_k(\mathbf{X})) \right|^2 = 0.$$

Proof of Proposition 4.3. Since $|\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) - 1| \leq 1$, one has

$$\left| \left(\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) - 1 \right) T(\mathbf{r}_k(\mathbf{X})) \right|^2 \leq T^2(\mathbf{r}_k(\mathbf{X})).$$

Consequently, by Lebesgue's dominated convergence theorem, to prove the proposition, it suffices to show that $W_{n,i}(\mathbf{X})$ tends to 1 almost surely. Now,

$$\begin{aligned} & \mathbb{P} \left(\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) \neq 1 \right) \\ &= \mathbb{P} \left(\sum_{i=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{X}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_{\ell}\}} = 0 \right) \\ &= \mathbb{P} \left(\sum_{i=1}^{\ell} \mathbf{1}_{\{\mathbf{X}_i \in \cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{X}) - \varepsilon_{\ell}, r_{k,m}(\mathbf{X}) + \varepsilon_{\ell}])\}} = 0 \right) \\ &= \int_{\mathbb{R}^d} \mathbb{P} \left(\forall i = 1, \dots, \ell, \mathbf{1}_{\{\mathbf{X}_i \in \cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{x}) - \varepsilon_{\ell}, r_{k,m}(\mathbf{x}) + \varepsilon_{\ell}])\}} = 0 \right) \mu(d\mathbf{x}) \\ &= \int_{\mathbb{R}^d} [1 - \mu(\cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{x}) - \varepsilon_{\ell}, r_{k,m}(\mathbf{x}) + \varepsilon_{\ell}]))]^{\ell} \mu(d\mathbf{x}). \end{aligned}$$

Denote by I a bounded interval. Then,

$$\begin{aligned} & \mathbb{P} \left(\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) \neq 1 \right) \\ & \leq \int_{\mathbb{R}^d} \exp(-\ell \mu(\cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{x}) - \varepsilon_{\ell}, r_{k,m}(\mathbf{x}) + \varepsilon_{\ell}])) \\ & \quad \times \mathbf{1}_{\{\mathbf{x} \in \cap_{m=1}^M r_{k,m}^{-1}(I)\}} \mu(d\mathbf{x}) + \mu \left(\bigcup_{m=1}^M r_{k,m}^{-1}(I^c) \right) \\ & \leq \max_{\mathbf{u}} \mathbf{u} e^{-\mathbf{u}} \int_{\mathbb{R}^d} \frac{\mathbf{1}_{\{\mathbf{x} \in \cap_{m=1}^M r_{k,m}^{-1}(I)\}}}{\ell \mu(\cap_{m=1}^M r_{k,m}^{-1}([r_{k,m}(\mathbf{x}) - \varepsilon_{\ell}, r_{k,m}(\mathbf{x}) + \varepsilon_{\ell}]))} \mu(d\mathbf{x}) \\ & \quad + \mu \left(\bigcup_{m=1}^M r_{k,m}^{-1}(I^c) \right). \end{aligned}$$

Using the same arguments as in the proof of Proposition 4.2, the probability $\mathbb{P} \left(\sum_{i=1}^{\ell} W_{n,i}(\mathbf{X}) \neq 1 \right)$ is bounded by $\frac{e^{-1}}{\ell} \left[\frac{|I|}{\varepsilon_{\ell}} \right]^M$. This bound vanishes as n tends to infinity since, by assumption, $\lim_{\ell \rightarrow \infty} \ell \varepsilon_{\ell}^M = \infty$. \square

4.3 Proof of Theorem 2.1

Choose $\mathbf{x} \in \mathbb{R}^d$. An easy calculation yields that

$$\begin{aligned}
& \mathbb{E}[|T_n(\mathbf{r}_k(\mathbf{x})) - T(\mathbf{r}_k(\mathbf{x}))|^2 | \mathbf{r}_k(\mathbf{X}_1), \dots, \mathbf{r}_k(\mathbf{X}_\ell), \mathcal{D}_k] \\
&= \mathbb{E} \left[|T_n(\mathbf{r}_k(\mathbf{x})) - \mathbb{E}[T_n(\mathbf{r}_k(\mathbf{x})) | \mathbf{r}_k(\mathbf{X}_1), \dots, \mathbf{r}_k(\mathbf{X}_\ell), \mathcal{D}_k]|^2 | \mathbf{r}_k(\mathbf{X}_1), \dots, \mathbf{r}_k(\mathbf{X}_\ell), \mathcal{D}_k \right] \\
&\quad + |\mathbb{E}[T_n(\mathbf{r}_k(\mathbf{x})) | \mathbf{r}_k(\mathbf{X}_1), \dots, \mathbf{r}_k(\mathbf{X}_\ell), \mathcal{D}_k] - T(\mathbf{r}_k(\mathbf{x}))|^2 \\
&:= E_1 + E_2.
\end{aligned} \tag{4.6}$$

On the one hand, we have

$$\begin{aligned}
E_1 &= \mathbb{E} \left[|T_n(\mathbf{r}_k(\mathbf{x})) - \mathbb{E}[T_n(\mathbf{r}_k(\mathbf{x})) | \mathbf{r}_k(\mathbf{X}_1), \dots, \mathbf{r}_k(\mathbf{X}_\ell), \mathcal{D}_k]|^2 | \mathbf{r}_k(\mathbf{X}_1), \dots, \mathbf{r}_k(\mathbf{X}_\ell), \mathcal{D}_k \right] \\
&= \mathbb{E} \left[\left| \sum_{i=1}^{\ell} W_{n,i}(\mathbf{x}) (Y_i - \mathbb{E}[Y_i | \mathbf{r}_k(\mathbf{X}_i)]) \right|^2 | \mathbf{r}_k(\mathbf{X}_1), \dots, \mathbf{r}_k(\mathbf{X}_\ell), \mathcal{D}_k \right].
\end{aligned}$$

Developing the square and noticing that $\mathbb{E}[Y_j | Y_i, \mathbf{r}_k(\mathbf{X}_1), \dots, \mathbf{r}_k(\mathbf{X}_\ell), \mathcal{D}_k] = \mathbb{E}[Y_j | \mathbf{r}_k(\mathbf{X}_j)]$, since Y_j is independent of Y_i and of the X_j 's with $j \neq i$, we have

$$\begin{aligned}
E_1 &= \mathbb{E} \left[\frac{\sum_{i=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_\ell\}} |Y_i - \mathbb{E}[Y_i | \mathbf{r}_k(\mathbf{X}_i)]|^2}{\left| \sum_{i=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_\ell\}} \right|^2} \right. \\
&\quad \left. | \mathbf{r}_k(\mathbf{X}_1), \dots, \mathbf{r}_k(\mathbf{X}_\ell), \mathcal{D}_k \right] \\
&= \sum_{i=1}^{\ell} \mathbb{V}(Y_i | \mathbf{r}_k(\mathbf{X}_i)) \frac{\mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_\ell\}}}{\left| \sum_{i=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_\ell\}} \right|^2}.
\end{aligned} \tag{4.7}$$

Thus,

$$E_1 \leq 4R^2 \frac{\mathbf{1}_{\left\{ \sum_{i=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_\ell\}} > 0 \right\}}}{\sum_{i=1}^{\ell} \mathbf{1}_{\cap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_\ell\}}}, \tag{4.8}$$

where $\mathbb{V}(Z)$ denotes the variance of a random variable Z . On the other hand, and recalling the notation Σ introduced in Section 3, we obtain for the second

term E_2 :

$$\begin{aligned}
E_2 &= \left| \mathbb{E}[T_n(\mathbf{r}_k(\mathbf{x})) | \mathbf{r}_k(\mathbf{X}_1), \dots, \mathbf{r}_k(\mathbf{X}_\ell), \mathcal{D}_k] - T(\mathbf{r}_k(\mathbf{x})) \right|^2 \\
&= \left| \sum_{i=1}^{\ell} W_{n,i}(\mathbf{x}) \mathbb{E}[Y_i | \mathbf{r}_k(\mathbf{X}_i)] - T(\mathbf{r}_k(\mathbf{x})) \right|^2 \mathbf{1}_{\{\Sigma > 0\}} + T^2(\mathbf{r}_k(\mathbf{x})) \mathbf{1}_{\{\Sigma = 0\}} \\
&\leq \frac{\sum_{i=1}^{\ell} \mathbf{1}_{\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_\ell\}} |\mathbb{E}[Y_i | \mathbf{r}_k(\mathbf{X}_i)] - T(\mathbf{r}_k(\mathbf{x}))|^2}{\sum_{j=1}^{\ell} \mathbf{1}_{\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_j)| \leq \varepsilon_\ell\}}} \mathbf{1}_{\{\Sigma > 0\}} \quad (4.9)
\end{aligned}$$

$$\begin{aligned}
&+ T^2(\mathbf{r}_k(\mathbf{x})) \mathbf{1}_{\{\Sigma = 0\}} \\
&\quad (\text{by Jensen's inequality}) \\
&= \frac{\sum_{i=1}^{\ell} \mathbf{1}_{\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_i)| \leq \varepsilon_\ell\}} |T(\mathbf{r}_k(\mathbf{X}_i)) - T(\mathbf{r}_k(\mathbf{x}))|^2}{\sum_{j=1}^{\ell} \mathbf{1}_{\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X}_j)| \leq \varepsilon_\ell\}}} \mathbf{1}_{\{\Sigma > 0\}} \quad (4.10)
\end{aligned}$$

$$\begin{aligned}
&+ T^2(\mathbf{r}_k(\mathbf{x})) \mathbf{1}_{\{\Sigma = 0\}} \\
&\leq L^2 \varepsilon_\ell^2 + T^2(\mathbf{r}_k(\mathbf{x})) \mathbf{1}_{\{\Sigma = 0\}}. \quad (4.11)
\end{aligned}$$

Now,

$$\mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 \leq \int_{\mathbb{R}^d} \mathbb{E}|(T_n(\mathbf{r}_k(\mathbf{x})) - T(\mathbf{r}_k(\mathbf{x})))|^2 \mu(d\mathbf{x}).$$

Then, using the decomposition (4.6) and the upper bounds (4.8) and (4.11),

$$\begin{aligned}
&\mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 \\
&\leq \int_{\mathbb{R}^d} \mathbb{E} \left[\frac{4R^2 \mathbf{1}_{\{\Sigma > 0\}}}{B} \right] \mu(d\mathbf{x}) + L^2 \varepsilon_\ell^2 + \int_{\mathbb{R}^d} \mathbb{E} [T^2(\mathbf{r}_k(\mathbf{x})) \mathbf{1}_{\{\Sigma = 0\}}] \mu(d\mathbf{x}) \\
&\leq \int_{\mathbb{R}^d} \mathbb{E} \left\{ \mathbb{E} \left[\frac{4R^2 \mathbf{1}_{\{\Sigma > 0\}}}{B} \middle| \mathcal{D}_k \right] \right\} \mu(d\mathbf{x}) + L^2 \varepsilon_\ell^2 \\
&\quad + \int_{\mathbb{R}^d} \mathbb{E} \left\{ \mathbb{E} [T^2(\mathbf{r}_k(\mathbf{x})) \mathbf{1}_{\{\Sigma = 0\}} | \mathcal{D}_k] \right\} \mu(d\mathbf{x}).
\end{aligned}$$

Thus, thanks to Lemma 4.1,

$$\begin{aligned}
&\mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 \\
&\leq \frac{8R^2}{(\ell + 1)} \int_{\mathbb{R}^d} \frac{1}{\mu(\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X})| \leq \varepsilon_\ell\})} \mu(d\mathbf{x}) + L^2 \varepsilon_\ell^2 \\
&\quad + \int_{\mathbb{R}^d} T^2(\mathbf{r}_k(\mathbf{x})) \left(1 - \mu\left(\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X})| \leq \varepsilon_\ell\}\right) \right)^\ell \mu(d\mathbf{x}).
\end{aligned}$$

Consequently,

$$\begin{aligned}
& \mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 \\
& \leq \frac{8R^2}{(\ell+1)} \int_{\mathbb{R}^d} \frac{1}{\mu(\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X})| \leq \varepsilon_\ell\})} \mu(d\mathbf{x}) + L^2 \varepsilon_\ell^2 \\
& \quad + \int_{\mathbb{R}^d} T^2(\mathbf{r}_k(\mathbf{x})) \exp \left(-\ell \mu \left(\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X})| \leq \varepsilon_\ell\} \right) \right) \mu(d\mathbf{x}) \\
& \leq \frac{8R^2}{(\ell+1)} \int_{\mathbb{R}^d} \frac{1}{\mu(\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X})| \leq \varepsilon_\ell\})} \mu(d\mathbf{x}) + L^2 \varepsilon_\ell^2 \\
& \quad + \left(\sup_{\mathbf{x} \in \mathbb{R}^d} T^2(\mathbf{r}_k(\mathbf{x})) \max_{\mathbf{u} \in \mathbb{R}^+} \mathbf{u} e^{-\mathbf{u}} \right. \\
& \quad \left. \times \int_{\mathbb{R}^d} \frac{1}{\ell \mu(\bigcap_{m=1}^M \{|r_{k,m}(\mathbf{x}) - r_{k,m}(\mathbf{X})| \leq \varepsilon_\ell\})} \mu(d\mathbf{x}) \right).
\end{aligned}$$

Introducing a bounded interval I as in the proof of [Proposition 2.2](#), we observe that the boundedness of the \mathbf{r}_k yields that

$$\mu \left(\bigcup_{m=1}^M r_{k,m}^{-1}(I^c) \right) = 0,$$

as soon as I is sufficiently large, independently of k . Then, proceeding as in the proof of [Proposition 2.2](#), we obtain

$$\begin{aligned}
& \mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 \\
& \leq 8R^2 \left\lceil \frac{|I|}{\varepsilon_\ell} \right\rceil^M \frac{1}{\ell+1} + L^2 \varepsilon_\ell^2 + R^2 \max_{\mathbf{u} \in \mathbb{R}^+} \mathbf{u} e^{-\mathbf{u}} \left\lceil \frac{|I|}{\varepsilon_\ell} \right\rceil^M \frac{1}{\ell} \\
& \leq C_1 \frac{R^2}{\ell \varepsilon_\ell^M} + L^2 \varepsilon_\ell^2,
\end{aligned}$$

for some positive constant C_1 , independent of k . Hence, for the choice $\varepsilon_\ell \propto \ell^{-\frac{1}{M+2}}$, we obtain

$$\mathbb{E}|T_n(\mathbf{r}_k(\mathbf{X})) - T(\mathbf{r}_k(\mathbf{X}))|^2 \leq C \ell^{-\frac{2}{M+2}},$$

for some positive constant C depending on L , R and independent of k , as desired.

Acknowledgements

The authors thank the Joint Editor and three anonymous referees for providing constructive and helpful remarks, thus greatly improving the paper.

References

- ALQUIER, P. and BIAU, G. (2013). Sparse Single-Index Model. *Journal of Machine Learning Research*, **14** 243–280.
- AUDIBERT, J.-Y. (2004). Aggregated estimators and empirical complexity for least square regression. *Annales de l’Institut Henri Poincaré : Probabilités et Statistiques*, **40** 685–736.
- BARAUD, Y., GIRAUD, C. and HUET, S. (2013). Estimator selection in the Gaussian setting. *Annales de l’Institut Henri Poincaré (B) Probabilités et Statistiques*.
- BIAU, G. (2012). Analysis of a random forests model. *Journal of Machine Learning Research*, **13** 1063–1095.
- BIRGÉ, L. (2006). Model selection via testing: An alternative to (penalized) maximum likelihood estimators. *Annales de l’Institut Henri Poincaré (B) Probabilités et Statistiques*, **42** 273–325.
- BUNEA, F. and NOBEL, A. (2008). Sequential procedures for aggregating arbitrary estimators of a conditional mean. *IEEE Transactions on Information Theory*, **54** 1725–1735.
- BUNEA, F., TSYBAKOV, A. B. and WEGKAMP, M. H. (2006). Aggregation and sparsity via ℓ_1 -penalized least squares. In *Proceedings of the 19th annual conference on Computational Learning Theory* (G. Lugosi and H. U. Simon, eds.), vol. 35. Springer-Verlag, 379–391.
- BUNEA, F., TSYBAKOV, A. B. and WEGKAMP, M. H. (2007a). Aggregation for gaussian regression. *The Annals of Statistics*, **35** 1674–1697.
- BUNEA, F., TSYBAKOV, A. B. and WEGKAMP, M. H. (2007b). Sparsity oracle inequalities for the Lasso. *Electronic Journal of Statistics*, **35** 169–194.
- CATONI, O. (2004). *Statistical Learning Theory and Stochastic Optimization*. École d’Été de Probabilités de Saint-Flour XXXI – 2001, Springer.
- CORTEZ, P., CERDEIRA, A., ALMEIDA, F., MATOS, T. and REIS, J. (2009). Modeling wine preferences by data mining from physicochemical properties. *Decision Support Systems*, **47** 547–553.
- CULE, E. (2012). *ridge: Ridge Regression with automatic selection of the penalty parameter*. R package version 2.1-2, URL <http://CRAN.R-project.org/package=ridge>.

- DALALYAN, A. S. and TSYBAKOV, A. B. (2008). Aggregation by exponential weighting, sharp PAC-Bayesian bounds and sparsity. *Machine Learning*, **72** 39–61.
- DEVROYE, L., GYÖRFI, L. and LUGOSI, G. (1996). *A probabilistic theory of pattern recognition*. Springer.
- GUEDJ, B. (2013). *COBRA: COmBined Regression Alternative*. R package version 0.99.4, URL <http://cran.r-project.org/web/packages/COBRA/index.html>.
- GUEDJ, B. and ALQUIER, P. (2013). PAC-Bayesian estimation and prediction in sparse additive models. *Electronic Journal of Statistics*, **7** 264–291.
- GYÖRFI, L., KOHLER, M., KRZYŻAK, A. and WALK, H. (2002). *A Distribution-Free Theory of Nonparametric Regression*. Springer Series in Statistics, Springer.
- HASTIE, T. and EFRON, B. (2012). *lars: Least Angle Regression, Lasso and Forward Stagewise*. R package version 1.1, URL <http://CRAN.R-project.org/package=lars>.
- JUDITSKY, A., NAZIN, A. V., TSYBAKOV, A. B. and VAYATIS, N. (2005). Recursive aggregation of estimators by the mirror descent method with averaging. *Problems of Information Transmission*, **41** 368–384.
- JUDITSKY, A. and NEMIROVSKI, A. (2000). Functional aggregation for nonparametric estimation. *The Annals of Statistics*, **28** 681–712.
- KNAUS, J. (2010). *snowfall: Easier cluster computing (based on snow)*. R package version 1.84, URL <http://CRAN.R-project.org/package=snowfall>.
- KOLTCHINSKII, V. (2009). Sparsity in penalized empirical risk minimization. *Annales de l’Institut Henri Poincaré (B) Probabilités et Statistiques*, **45** 7–57.
- LI, S. (2013). *FNN: Fast Nearest Neighbor Search Algorithms and Applications*. R package version 1.1, URL <http://CRAN.R-project.org/package=FNN>.
- LIAW, A. and WIENER, M. (2002). Classification and regression by random-forest. *R News*, **2** 18–22. URL <http://CRAN.R-project.org/doc/Rnews/>.
- MALLEY, J. D., KRUPPA, J., DASGUPTA, A., MALLEY, K. G. and ZIEGLER, A. (2012). Probability machines: Consistent probability estimation using nonparametric learning machines. *Methods of Information in Medicine*, **51** 74–81.

- MEIER, L., VAN DE GEER, S. A. and BÜHLMANN, P. (2009). High-dimensional additive modeling. *The Annals of Statistics*, **37** 3779–3821.
- MOJIRSHEIBANI, M. (1999). Combining classifiers via discretization. *Journal of the American Statistical Association*, **94** 600–609.
- NEMIROVSKI, A. (2000). *Topics in Non-Parametric Statistics*. Springer.
- PEKALSKA, E. and DUIN, R. P. W. (2005). *The Dissimilarity Representation For Pattern Recognition: Foundations and Applications*, vol. 64 of *Machine Perception and Artificial Intelligence*. World Scientific.
- POLLEY, E. C. and VAN DER LAAN, M. J. (2010). Super learner in prediction. Tech. rep., UC Berkeley.
- POLLEY, E. C. and VAN DER LAAN, M. J. (2012). *SuperLearner: Super Learner Prediction*. R package version 2.0-9, URL <http://CRAN.R-project.org/package=SuperLearner>.
- R CORE TEAM (2013). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>.
- RIPLEY, B. (2012). *tree: Classification and regression trees*. R package version 1.0-32, URL <http://CRAN.R-project.org/package=tree>.
- TSYBAKOV, A. B. (2003). Optimal rates of aggregation. In *Computational Learning Theory and Kernel Machines* (B. Schölkopf and M. K. Warmuth, eds.). Lecture Notes in Computer Science, Springer-Verlag, Berlin Heidelberg, Springer, Heidelberg, 303–313.
- VAN DE GEER, S. A. (2008). High-dimensional generalized linear models and the Lasso. *The Annals of Statistics*, **36** 614–645.
- VAN DER LAAN, M. J., POLLEY, E. C. and HUBBARD, A. E. (2007). Super learner. *Statistical Applications in Genetics and Molecular Biology*, **6**.
- WEGKAMP, M. H. (2003). Model selection in nonparametric regression. *The Annals of Statistics*, **31** 252–273.
- YANG, Y. (2000). Combining different procedures for adaptive regression. *Journal of Multivariate Analysis*, **74** 135–161.
- YANG, Y. (2001). Adaptive regression by mixing. *Journal of the American Statistical Association*, **96** 574–588.
- YANG, Y. (2004). Aggregating regression procedures to improve performance. *Bernoulli*, **10** 25–47.

- YEH, I.-C. (1998). Modeling of strength of high performance concrete using artificial neural networks. *Cement and Concrete Research*, **28** 1797–1808.
- YEH, I.-C. (2007). Modeling slump flow of concrete using second-order regressions and artificial neural networks. *Cement and Concrete Composites*, **29** 474–480.