# ON ADAPTIVE ESTIMATION OF NONPARAMETRIC FUNCTIONALS

By Rajarshi Mukherjee[*], Eric Tchetgen Tchetgen[†], and James Robins[‡]

*Abstract* We provide general adaptive upper bounds for estimating nonparametric functionals based on second order U-statistics arising from finite dimensional approximation of the infinite dimensional models using projection type kernels. An accompanying general adaptive lower bound tool is provided yielding bounds on chi-square divergence between mixtures of product measures. We then provide examples of functionals for which the theory produces rate optimally matching adaptive upper and lower bound.

**1. Introduction.** Estimation of functionals of data generating distribution has always been of central interest in statistics. In nonparametric statistics, where data generating distributions are parametrized by functions in infinite dimensional spaces, there exists a comprehensive literature addressing such questions. In particular, a large body of research has been devoted to explore minimax estimation of linear and quadratic functionals in density and white noise models. We do not attempt to survey the extensive literature in this area. However, the interested reader can find a comprehensive snapshot of the literature in Hall and Marron (1987), Bickel and Ritov (1988), Donoho, Liu and MacGibbon (1990), Donoho and Nussbaum (1990), Fan (1991), Kerkyacharian and Picard (1996), Laurent (1996), Cai and Low (2003), Cai and Low (2004), Cai and Low (2005a), and other references therein. Although the question of more general nonparametric functionals has received relatively less attention, some fundamental insights regarding estimation of non-linear integral functionals in density and white noise models can be found in Ibragimov and Has' minskii (2013), Kerkyacharian and Picard (1996), Nemirovski (2000), and references therein.

A general feature of the results obtained while estimating "smooth" nonparametric functionals is an elbow effect in the rate of estimation based on the smoothness of the underlying function classes. For example while estimating quadratic functionals in a $d$ dimensional density model, $\sqrt{n}$-efficient estimation can be achieved as soon as Hölder exponent $\beta$ of the underlying density exceeds $\frac{d}{4}$, whereas the optimal rate of estimation is $n^{-\frac{4\beta}{4\beta+d}}$ (in root mean squared sense) for $\beta < \frac{d}{4}$. A similar elbow in the rate of estimation exists for estimation of non-linear integral functionals as well. For density model this was demonstrated by Birgé and Massart (1995), Kerkyacharian and Picard (1996). For signal or white noise model, the problem of general integrated non-linear functionals was studied by Nemirovski (2000), but mostly in the $\sqrt{n}$- regime. However, for more complex nonparametric models, the approach for constructing minimax optimal procedures for general non-linear functionals in non-$\sqrt{n}$ regimes has been rather case specific. Motivated by this, in recent times, Robins et al. (2008, 2016), Mukherjee, Newey and Robins (2017) have developed a theory of inference for nonlinear functionals in parametric, semi-parametric, and non-parametric models based on higher order influence functions.

Most minimax rate optimal estimators proposed in the literature, however, depend explicitly on the knowledge of the smoothness indices. Thus, it becomes of interest to understand the question of adaptive estimation i.e. the construction and analysis of estimators without prior knowledge of

the smoothness. The question of adaptation of linear and quadratic functionals has been studied in detail in the context of density, white noise, and nonparametric additive Gaussian noise regression models (Low (1992), Efromovich and Low (1994), Efromovich and Low (1996), Tribouley (2000), Efromovich and Samarov (2000), Klemela and Tsybakov (2001), Laurent and Massart (2000), Cai and Low (2005b), Cai and Low (2006), Giné and Nickl (2008)). However, adaptive estimation of more general non-linear functionals in more complex nonparametric models have received much less attention. This paper is motivated by taking a step in that direction.

In particular, we suppose we observe i.i.d copies of a random vector $O = (\mathbf{W}; \mathbf{X}) \in \mathbf{R}^{m+d}$ with unknown distribution $P$ on each of $n$ study subjects. The variable $\mathbf{X}$ represents a random vector of baseline covariates such as age, height, weight, etc. Throughout $\mathbf{X}$ is assumed to have compact support and a density with respect to the Lebesgue measure in $\mathbb{R}^d$. The variable $\mathbf{W} \in \mathbb{R}^m$ can be thought of as a combination of outcome and treatment variables in some of our examples. In the above set up, we are interested in estimating certain "smooth" functionals $\phi(P)$ in the sense that under finite dimensional parametric submodels, they admit derivatives which can be represented as inner products of first order influence functions with score functions (Bickel et al., 1993). For some classic examples of these functionals, we provide matching upper and lower bounds on the rate of adaptive minimax estimation over a varying class of smoothness of the underlying functions, provided the marginal design density of $\mathbf{X}$ is sufficiently smooth.

The contributions of this paper are as follows. Extending the theory from density estimation and Gaussian white noise models, we provide a step towards adaptation theory for non-linear functionals in more complex nonparametric models in the non-$\sqrt{n}$ regime. The crux of our arguments relies on the observation that when the non-adaptive minimax estimators can be written as a sum of empirical mean type statistics and $2^{\text{nd}}$-order U-statistics, then one can provide a unified theory of selecting the "best" data driven estimator using Lepski type arguments (Lepski, 1991, 1992). Indeed, under certain assumptions on the data generating mechanism $P$, the non-adaptive minimax estimators have the desired structure for a large class of problems (Robins et al., 2008). This enables us to produce a class of examples where a single method helps provide a desired answer. In order to prove a lower bound for the rate of adaptation, we provide a general tool for bounding the chi-square divergence between two mixtures of suitable product measures. This extends the results in Birgé and Massart (1995); Robins et al. (2009), where similar results were obtained for the Hellinger distance. Our results are provided in the low regularity regime, i.e. when it is not possible to achieve $\sqrt{n}$-efficient estimator in an asymptotically minimax sense. This typically happens when the "average smoothness" of the function classes in consideration is below $\frac{d}{4}$. Discussions on obtaining a corresponding $\sqrt{n}$-efficient estimator for regularity above $\frac{d}{4}$ is provided in Section 4.

The rest of the paper is organized as follows. In Section 2 we provide the main results of the paper in a general form. Section 3 is devoted for applications of the main results in specific examples. A discussion on some issues left unanswered is provided in Section 4. In Section 5 we provide a brief discussion on some basic wavelet and function space theory and notations, which we use extensively in our jargon. Finally Section 6 and Appendices A and B are devoted for the proof of the theorems and collecting useful technical lemmas.

1.1. *Notation.* For data arising from underlying distribution $P$ we denote by $\mathbb{P}_P$ and $\mathrm{E}_P$ the probability of an event and expectation under $P$ receptively. For any positive integer $m \geq 2^d$, let $j(m)$ denote the largest integer $j$ such that $2^{jd} \leq m$ i.e. $j = \lfloor \frac{1}{d} \frac{\log m}{\log 2} \rfloor$ and $2^{j(m)d} \geq m/2^d$. For a two variable function $h(O_1, O_2)$ let $S(h(O_1, O_2)) = \frac{1}{2}[h(O_1, O_2) + h(O_2, O_1)]$ be the symmetrization of $h$. The results in this paper are mostly asymptotic (in $n$) in nature and thus requires some standard asymptotic notations. If $a_n$ and $b_n$ are two sequences of real numbers then $a_n \gg b_n$ (and $a_n \ll b_n$) implies that $a_n/b_n \to \infty$ (and $a_n/b_n \to 0$) as $n \to \infty$, respectively. Similarly $a_n \gtrsim b_n$ (and $a_n \lesssim b_n$)

implies that $\liminf a_n/b_n = C$ for some $C \in (0, \infty]$ (and $\limsup a_n/b_n = C$ for some $C \in [0, \infty)$). Alternatively, $a_n = o(b_n)$ will also imply $a_n \ll b_n$ and $a_n = O(b_n)$ will imply that $\limsup a_n/b_n = C$ for some $C \in [0, \infty)$). Finally we comment briefly on the various constants appearing throughout the text and proofs. Given that our primary results concern convergence rates of various estimators, we will not emphasize the role of constants throughout and rely on fairly generic notation for such constants. In particular, for any fixed tuple $v$ of real numbers, $C(v)$ will denote a positive real number which depends on elements of $v$ only. Finally for any linear subspace $L \subseteq L_2[0,1]^d$, let $\Pi(h|L)$ denote the orthogonal projection of $h$ onto $L$ under the Lebesgue measure. Also, for a function defined on $[0,1]^d$, for $1 \le q < \infty$ we let $\|h\|_q := (\int_{[0,1]^d} |h(\mathbf{x})|^q d\mathbf{x})^{1/q}$ denote the $L_q$ semi-norm of $h$, $\|h\|_\infty := \sup_{\mathbf{x} \in [0,1]^d} |h(\mathbf{x})|$ the $L_\infty$ semi-norm of $h$. We say $h \in L_q[0,1]^d$ for $q \in [1, \infty]$ if $\|h\|_q < \infty$. Typical functions arising in this paper will be considered to have memberships in certain Hölder balls $H(\beta, M)$ (see Section 5 for details). This will imply that the functions are uniformly bounded by a number depending on $M$. However, to make the dependence of our results on the uniform upper bound of functions more clear, we will typically assume a bound $B_U$ over the function classes, and for the sake of compactness will avoid the notational dependence of $B_U$ on $M$.

**2. Main Results.** We divide the main results of the paper into three main parts. First we discuss a general recipe for producing a "best" estimator from a sequence of estimators based on second order U-statistics constructed from compactly supported wavelet based projection kernels (defined in Section 5). Next we provide a general tool for bounding chi-square divergence between mixtures of product measures. This serves as a basis of using a version of constrained risk inequality (Cai and Low, 2011) for producing matching adaptive lower bounds in context of estimation of non-linear functionals considered in this paper. Finally we provide estimators of design density as well as regression functions in $L_\infty$ norm which adapt over Hölder type smoothness classes (defined in Section 5).

2.1. *Upper Bound.*
Consider i.i.d data $O_i = (\mathbf{W}_i, \mathbf{X}_i) \sim P$, $\mathbf{W}_i \in \mathbb{R}^m$, $\mathbf{X}_i \in [0,1]^d$, $i = 1, \ldots, n$ and a real valued functional of interest $\phi(P)$. Given this sample of size $n \ge 1$, consider further, a sequence of estimators $\{\hat{\phi}_{n,j}\}_{j \ge 1}$ of $\phi(P)$ defined as follows:

$$\hat{\phi}_{n,j} = \frac{1}{n} \sum_{i=1}^{n} L_1(O_i) - \frac{1}{n(n-1)} \sum_{i_1 \ne i_2} S\left(L_{2l}(O_{i_1}) K_{V_j}(\mathbf{X}_{i_1}, \mathbf{X}_{i_2}) L_{2r}(O_{i_2})\right)$$

where $K_{V_j}(\mathbf{X}_{i_1}, \mathbf{X}_{i_2})$ is a resolution $2^j$ wavelet projection kernel defined in Section 5 and $L_1, L_{2l}, L_{2r}$ are measurable functions such that $\forall O$ one has

$$\max\{|L_1(O)|, \ |L_{2l}(O)|, \ |L_{2r}(O)|\} \le B$$

for a known constant $B$. Also assume that $|g(\mathbf{x})| \le B \ \forall \mathbf{x}$, $g$ being the marginal density of $\mathbf{X}$ with respect to Lebesgue measure. Such a sequence of estimators can be though about as a bias corrected version of a usual first order estimator arising from standard first order influence function theory for "smooth" functionals $\phi(P)$ (Bickel et al., 1993; Van der Vaart, 2000). In particular, the linear empirical mean type term $\frac{1}{n} \sum_{i=1}^{n} L_1(O_i)$ typically derives from the classical influence function of $\phi(P)$ and the U-statistics quadratic term $\frac{1}{n(n-1)} \sum_{i_1 \ne i_2} S\left(L_{2l}(O_{i_1}) K_{V_j}(\mathbf{X}_{i_1}, \mathbf{X}_{i_2}) L_{2r}(O_{i_2})\right)$ corrects for higher order bias terms. While, specific examples in Section 3 will make the structure of these sequence of estimators more clear, the interested reader will be able to find more detailed theory in Robins et al. (2008, 2016).

The quality of such sequence of estimators will be judged against models for the data generating mechanism $P$. To this end, assume $P \in \mathcal{P}_\theta$ where $\mathcal{P}_\theta$ is a class of data generating distributions indexed by $\theta$ which in turn can vary over an index set $\Theta$. The choices of such a $\Theta$ will be clear from our specific examples in Section 3, and will typically corresponding to smoothness indices of various infinite dimensional functions parametrizing the data generating mechanisms. Further assume that there exists positive real values functions $f_1$ and $f_2$ defined on $\Theta$ such that the sequence of estimators $\{\hat{\phi}_{n,j}\}_{j \geq 1}$ satisfies the following bounds with known constants $C_1 > 0$ and $C_2 > 0$ whenever $n \leq 2^{jd} \leq n^2$.

Property (A): Bias Bound

$$\sup_{P \in \mathcal{P}_\theta} |\mathrm{E}_P \left( \hat{\phi}_{n,j} - \phi(P) \right)| \leq C_1 \left( 2^{-2jd\frac{f_1(\theta)}{d}} + n^{-f_2(\theta)} \right).$$

Property (B): Variance Bound

$$\sup_{P \in \mathcal{P}_\theta} \mathrm{E}_P \left( \hat{\phi}_{n,j} - \mathrm{E}_P \left( \hat{\phi}_{n,j} \right) \right)^2 \leq C_2 \frac{2^{jd}}{n^2}.$$

Given properties (A) and (B), we employ a Lepski type argument to choose an "optimal" estimator from the collection of $\{\hat{\phi}_{n,j}\}_{j \geq 1}$'s. To this end, as in Efromovich and Low (1996), for a given choice of $c > 1$, let $N$ be the largest integer such that $c^{N-1} \leq n^{1 - \frac{2}{\log \log n}}$. Denote $k(j) = 2^{jd}$, which, according to the definition in Section 1.1 implies that $k(j(m)) \leq m$. For $l = 0, \ldots, N-1$ let $\beta_l$ be the solution of $k(j_l) = n^{\frac{2}{1+4\beta_l/d}}$, where $j_l = j(c^l n)$ i.e. $k(j_l) = 2^{j(c^l n)d}$. Note that, $\frac{k(j_N)}{n^2} = \frac{2^{j(c^{N-1}n)d}}{n^2} \leq \frac{c^{N-1}n}{n^2} = \frac{1}{n^{\frac{2}{\log \log n}}} = o(1)$. By our choice of discretization, $k(j_0) \leq k(j_1) \leq \ldots \leq k(j_{N-1})$. Also, there exists constants $c_1, c_2$ such that for all $0 \leq l \leq l' \leq N-1$ one has $\frac{\beta_l}{d} - \frac{\beta_{l'}}{d} \in \left[ c_1 \frac{l'-l}{\log n}, c_2 \frac{l'-l}{\log n} \right]$. For $l = 0, \ldots, N-1$ let $k_*(j_l) = \frac{k(j_l)}{(\log n)^{\frac{1}{1+4\beta_l/d}}} = \left( \frac{n^2}{\log n} \right)^{\frac{1}{1+4\beta_l/d}}$ and $R(k_*(j_l)) = \frac{k_*(j_l)}{n^2} = n^{-\frac{8\beta_l/d}{1+4\beta_l/d}} (\log n)^{-\frac{1}{1+4\beta_l/d}}$. This implies that for $l > l'$ there exists constant $c' \geq 0$ such that

$$\frac{k_*(j_l)}{k_*(j_{l'})} \geq \left( \frac{n^2}{\log n} \right)^{\frac{4\beta_{l'}/d - 4\beta_l/d}{(1+4\beta_{l'}/d)(1+4\beta_l/d)}} \geq e^{c'} \geq 1.$$

Finally letting $s^*(n)$ be the smallest $l$ such that $k_*(j_l) \geq n$, define

$$\hat{l} := \min \left\{ \begin{array}{c} l \colon \left( \hat{\phi}_{n,j(k_*(j_l))} - \hat{\phi}_{n,j(k_*(j_{l'}))} \right)^2 \leq C_{\mathrm{opt}}^2 R(k_*(j_{l'})) \log n \\ \forall l' \geq l, s^* \leq l \leq N-1 \end{array} \right\} \tag{2.1}$$

where $C_{\mathrm{opt}}$ is a deterministic constant to be specified later.

With the notations and definitions as above we now have the following theorem which is the main result in the direction of adaptive upper bound in this paper.

THEOREM 2.1. *Assume $\beta < \frac{d}{4}$. Then there exists a positive $C$ depending on $B, C_1, C_2$, and choice of wavelet bases $\psi_{0,0}^0, \psi_{0,0}^1$ (defined in Section 5) such that*

$$\sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathrm{E}_P \left( \hat{\phi}_{n,j(k^*(\hat{l}))} - \phi(P) \right)^2 \leq C \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{8\beta}{4\beta+d}}.$$

A few remarks are in order regarding the implications of Theorem 2.1. In particular, provided one has a knowledge of a data generating $\theta$ and therefore of $f_1(\theta)$ and $f_2(\theta)$, one can use the bias and variance properties to achieve an optimal trade-off and subsequently obtain mean squared error in estimating $\phi(P)$ over $P \in \mathcal{P}_\theta$ which scales as $n^{-\frac{8\beta}{8\beta+d}}$. Theorem 2.1 demonstrates a logarithmic price payed by the estimator $\hat{\phi}_{n,j\left(k^*(\hat{i})\right)}$ in terms of estimating $\phi(P)$ over a class of data generating mechanisms $\left\{\mathcal{P}_\theta \colon f_1(\theta) = \beta, f_2(\theta) > \frac{4\beta}{4\beta+d}, \beta < \frac{d}{4}\right\}$. As will be demonstrated in Section 3, the term $f_1(\theta) = \beta$ usually drives the minimax rate of estimation whereas $f_2(\theta) > \frac{4\beta}{4\beta+d}$ is a regularity condition which typically will relate to marginal density of covariates in observational studies. Moreover, in our examples, the range of $\beta < d/4$ will be necessary to guarantee the non-existence of $\sqrt{n}$-consistent estimators of $\phi(P)$ over $P \in \mathcal{P}_\theta$ in a minimax sense.

Finally, it therefore remain to be explored whether this logarithmic price payed in Theorem 2.1 is indeed necessary. Using a chi-square divergence inequality developed in the next section along with a suitable version of constrained risk inequality (see Section B) we shall argue that the logarithmic price of Theorem 2.1 is indeed necessary for a class of examples naturally arising in many observational studies.

### 2.2. *Lower Bound.*

We provide a bound on the chi-square divergence between two mixture of product measures, and thereby extending results of Birgé and Massart (1995) and Robins et al. (2009). Since both Birgé and Massart (1995) and Robins et al. (2009) considered demonstrating bounds on the Hellinger divergence between mixtures of product measures, they do not automatically provide bounds on the corresponding chi-square divergences. However, such bounds are essential to explore adaptive lower bound in a mean squared error sense. To this end, as in Robins et al. (2009), let $O_1, \ldots, O_n$ be a random sample from a density $p$ with respect to measure $\mu$ on a sample space $(\chi, \mathcal{A})$. For $k \in \mathbb{N}$, let $\chi = \cup_{j=1}^k \chi_j$ be a measurable partition of the sample space. Given a vector $\lambda = (\lambda_1, \ldots, \lambda_k)$ in some product measurable space $\Lambda = \Lambda_1 \times \cdots \times \Lambda_k$ let $P_\lambda$ and $Q_\lambda$ be probability measures on $\chi$ such that

- $P_\lambda(\chi_j) = Q_\lambda(\chi_j) = p_j$ for every $\lambda$ and some $(p_1, \ldots, p_k)$ in the $k$-dimensional simplex.
- The restrictions $P_\lambda$ and $Q_\lambda$ to $\chi_j$ depends on the $j^{\text{th}}$ coordinate $\lambda_j$ of $\lambda = (\lambda_1, \ldots, \lambda_k)$ only.

For $p_\lambda$ and $q_\lambda$ densities of the measures $P_\lambda$ and $Q_\lambda$ respectively that are jointly measurable in the parameters $\lambda$ and the observations, and $\pi$ a probability measure on $\Lambda$, define $p = \int p_\lambda d\pi(\lambda)$ and $q_\lambda = \int q_\lambda d(\pi(\lambda))$, and set

$$a = \max_j \sup_\lambda \int_{\chi_j} \frac{(p_\lambda - p)^2}{p_\lambda} \frac{d\mu}{p_j},$$

$$b = \max_j \sup_\lambda \int_{\chi_j} \frac{(p_\lambda - p)^2}{p_\lambda} \frac{d\mu}{p_j},$$

$$\widetilde{c} = \max_j \sup_\lambda \int_{\chi_j} \frac{p^2}{p_\lambda} \frac{d\mu}{p_j},$$

$$\delta = \max_j \sup_\lambda \int_{\chi_j} \frac{(q-p)^2}{p_\lambda} \frac{d\mu}{p_j}.$$

With the notations and definitions as above we now have the following theorem which is the main result in the direction of adaptive lower bound of this paper.

THEOREM 2.2. *Suppose that $np_j(1 \vee a \vee b \vee \widetilde{c}) \leq A$ for all $j$ and for all $\lambda$, $\underline{B} \leq p_\lambda \leq \overline{B}$ for positive constants $A, \underline{B}, \overline{B}$. Then there exists a $C > 0$ that depends only on $A, \underline{B}, \overline{B}$, such that, for any product probability measure $\pi = \pi_1 \otimes \cdots \otimes \pi_k$, one has*

$$\chi^2 \left( \int P_\lambda d(\pi(\lambda)), \int Q_\lambda d(\pi(\lambda)) \right) \leq e^{Cn^2 (\max_j p_j)(b^2 + ab) + Cn\delta} - 1,$$

*where $\chi^2(\nu_1, \nu_2) = \int \left( \frac{d\nu_2}{d\nu_1} - 1 \right)^2 d\nu_1$ is the chi-square divergence between two probability measures $\nu_1$ and $\nu_2$ with $\nu_2 \ll \nu_1$.*

### 2.3. $L_\infty$-*Adaptive Estimation of Density and Regression Functions*.

We provide adaptive estimator of regression function in $L_\infty$ using Lepski type arguments (Lepski, 1992). Consider i.i.d data on $O_i = (W_i, \mathbf{X}_i) \sim P$ for a scalar variable $W$ such that $|W| \leq B_U$ and $E_P(W|\mathbf{X}) = f(\mathbf{X})$ almost surely, and $\mathbf{X} \in [0,1]^d$ has density $g$ such that $0 < B_L \leq g(\mathbf{x}) \leq B_U < \infty$ for all $\mathbf{x} \in [0,1]^d$. Although to be precise, we should put subscripts $P$ to $f, g$, we omit this since the context of their use is clear. We assume Hölder type smoothness (defined in Section 5) on $f, g$ and let $\mathcal{P}(\beta, \gamma) = \{P : (f, g) \in H(\beta, M) \times H(\gamma, M), |f(\mathbf{x})| \leq B_U, B_L \leq g(\mathbf{x}) \leq B_U \; \forall \; \mathbf{x} \in [0,1]^d\}$ denote classes of data generating mechanisms indexed by the smoothness indices. Then we have the following theorem, which considers adaptive estimation of $f, g$ in $L_\infty$ norm over $(\beta, \gamma) \in (\beta_{\min}, \beta_{\max}) \times (\gamma_{\min}, \gamma_{\max})$, for given positive real numbers $\beta_{\min}, \beta_{\max}, \gamma_{\min}, \gamma_{\max}$.

THEOREM 2.3. *If $\gamma_{\min} > \beta_{\max}$, then there exists an $\hat{f}$ and $\hat{g}$ depending on $M, B_L, B_U, \beta_{\min}, \beta_{\max}, \gamma_{\min}, \gamma_{\max}$, and choice of wavelet bases $\psi^0_{0,0}, \psi^1_{0,0}$ (defined in Section 5) such that the following hold for every $(\beta, \gamma) \in (\beta_{\min}, \beta_{\max}) \times (\gamma_{\min}, \gamma_{\max})$ with a large enough $C > 0$ depending possibly on $M, B_L, B_U$, and $\gamma_{\max}$.*

$$\sup_{P \in \mathcal{P}(\beta, \gamma)} E_P \|\hat{f} - f\|_\infty \leq (C)^{\frac{d}{2\beta + d}} \left( \frac{n}{\log n} \right)^{-\frac{\beta}{2\beta + d}},$$

$$E_P \|\hat{g} - g\|_\infty \leq (C)^{\frac{d}{2\gamma + d}} \left( \frac{n}{\log n} \right)^{-\frac{\gamma}{2\gamma + d}},$$

$$\sup_{P \in \mathcal{P}(\beta, \gamma)} \mathbb{P}_P \left( \hat{f} \notin H(\beta, C) \right) \leq \frac{1}{n^2},$$

$$\sup_{P \in \mathcal{P}(\beta, \gamma)} \mathbb{P}_P \left( \hat{g} \notin H(\gamma, C) \right) \leq \frac{1}{n^2},$$

$$|\hat{f}(\mathbf{x})| \leq 2B_U \text{ and } B_L/2 \leq \hat{g}(\mathbf{x}) \leq 2B_U \quad \forall \mathbf{x} \in [0,1]^d.$$

REMARK 2.4. *A close look at the proof of Theorem 2.3, shows that the proof continues to hold for $\beta_{\min} = 0$. Moreover, although we did not keep track of our constants, the purpose of keeping them in the form above is to show that the multiplicative constants is not arbitrarily big when $\beta$ is large.*

Theorems of the flavor of Theorem 2.3 is not uncommon in literature (see (Giné and Nickl, 2015, Chapter 8) for details). In particular, results of the kind stating that $\hat{g} \in H(\gamma, C)$ with high probability uniformly over $\mathcal{P}(\beta, \gamma)$ for a suitably large constant $C$ is often very easy to demonstrate. However, our proof shows that a suitably bounded estimator $\hat{g}$, which adapts over smoothness and satisfies $\hat{g} \in H(\gamma, C)$ with probability larger than $1 - \frac{1}{n^\kappa}$ uniformly over $\mathcal{P}(\beta, \gamma)$, for any $\kappa > 0$

and correspondingly large enough $C$. This result in turn turns out to be crucial for the purpose of controlling suitable bias terms in our functional estimation problems. Additionally, the results concerning $\hat{f}$ are relatively less common in an unknown design density setting. Indeed, adaptive estimation of regression function with random design over Besov type smoothness classes has been obtained by model selection type techniques by Baraud (2002) for the case of Gaussian errors. Our results in contrast hold for any regression model with bounded outcomes and compactly supported covariates having suitable marginal design density. Before proceeding we note that the results of this paper can be extended to include the whole class of doubly robust functionals considered in Robins et al. (2008). However we only provide specific examples here to demonstrate the clear necessity to pay a sharp poly-logarithmic penalty for adaptation in low regularity regimes.

**3. Examples.** In this section we discuss applications of Theorem 2.1 and Theorem 2.2 in producing rate optimal adaptive estimators of certain nonparametric functionals commonly arising in statistical literature. The proof of the results in this section can be found in Mukherjee, Tchetgen Tchetgen and Rob (2017).

3.1. *Average Treatment Effect.* In this subsection, we consider estimating the "treatment effect" of a treatment on an outcome in presence of multi-dimensional confounding variables (Crump et al., 2009; Robins, Mark and Newey, 1992). To be more specific, we consider a binary treatment $A$ and response $Y$ and $d$-dimensional covariate vector $\mathbf{X}$, and let $\tau$ be the variance weighted average treatment effect defined as

$$\tau := \mathrm{E}\left(\frac{Var(A|\mathbf{X})c(\mathbf{X})}{\mathrm{E}(Var(A|\mathbf{X}))}\right) = \frac{\mathrm{E}(cov(Y, A|\mathbf{X}))}{\mathrm{E}(Var(A|\mathbf{X}))}$$

Above

$$c(\mathbf{x}) = \mathrm{E}(Y|A = 1, \mathbf{X} = \mathbf{x}) - \mathrm{E}(Y|A = 0, \mathbf{X} = \mathbf{x}). \tag{3.1}$$

and under the assumption of no unmeasured confounding, $c(\mathbf{x})$ is often referred to as the average treatment effect among subjects with covariate value $\mathbf{X} = \mathbf{x}$. The reason of referring to $\tau$ as the average treatment effect can e further understood by considering semi-parametrically constrained model

$$c(\mathbf{x}) = \phi^* \text{ for all } \mathbf{x}, \tag{3.2}$$

or specifically the model

$$\mathrm{E}(Y|A, \mathbf{X}) = \phi^* A + b(\mathbf{X}).$$

It is clear that under (3.2) it turns out that , $\tau$ equals $\phi^*$. Moreover, the inference on $\tau$ is closely related to the estimation $\mathrm{E}(Cov(Y, A|\mathbf{X}))$ (Robins et al., 2008). Specifically, point and interval estimator for $\tau$ can be constructed from point and interval estimator of $\mathrm{E}(Cov(Y, A|X))$. To be more specific, for any fixed $\tau^* \in \mathbb{R}$, one can define $Y^*(\tau^*) = Y - \tau^* A$ and consider

$$\phi(\tau^*) = \mathrm{E}((Y^*(\tau^*) - \mathrm{E}(Y^*(\tau^*)|\mathbf{X}))(A - \mathrm{E}(A|\mathbf{X}))) = \mathrm{E}(cov(Y^*(\tau^*), A|\mathbf{X})).$$

it is easy to check that $\tau$ is the unique solution of $\phi(\tau^*) = 0$. Consequently, if we can construct estimator $\hat{\phi}(\tau^*)$ of $\phi(\tau^*)$, then $\hat{\tau}$ satisfying $\psi(\hat{\tau}) = 0$ is an estimator of $\tau$ with desirable properties. Moreover, $(1 - \alpha)$ confidence set for $\tau$ can be constructed as the set of values of $\tau^*$ for

which $(1 - \alpha)$ interval estimator of $\phi(\tau^*)$ contains the value 0. Finally, since $\mathrm{E}(Cov(Y, A|X)) = \mathrm{E}(\mathrm{E}(Y|X)\mathrm{E}(A|X)) - \mathrm{E}(AY)$, and $\mathrm{E}(AY)$ is estimable easily at a parametric rate, the crucial part of the problem hinges on the estimation of $\mathrm{E}(\mathrm{E}(Y|X)\mathrm{E}(A|X))$.

Henceforth, for the rest of the section, we assume that we observe $n$ iid copies of $O = (Y, A, \mathbf{X}) \sim P$ and we want to estimate $\phi(P) = \mathrm{E}_P(Cov_P(Y, A|\mathbf{X}))$. We assume that the marginal distribution of $\mathbf{X}$ has a density with respect to Lebesgue measure on $\mathbb{R}^d$ that has a compact support, which we assume to be $[0, 1]^d$ and let $g$ be the marginal density of $\mathbf{X}$ (i.e. $\mathrm{E}_P(h(\mathbf{X})) = \int_{[0,1]^d} h(\mathbf{x})g(\mathbf{x})d\mathbf{x}$ for all $P$-integrable function $h$), $a(\mathbf{X}) := \mathrm{E}_P(A|\mathbf{X})$, $b(\mathbf{X}) := \mathrm{E}_P(Y|\mathbf{X})$, and $c(\mathbf{X}) = \mathrm{E}_P(Y|A = 1, \mathbf{X}) - \mathrm{E}_P(Y|A = 0, \mathbf{X})$. Although to be precise, we should put subscripts $P$ to $a, b, g, c$, we omit this since the context of their use is clear. Let $\Theta := \{\theta = (\alpha, \beta, \gamma): \frac{\alpha+\beta}{2} < \frac{d}{4}, \gamma_{\max} \geq \gamma > \gamma_{\min} \geq 2(1 + \epsilon)\max\{\alpha, \beta\}\}$ for some fixed $\epsilon > 0$, and let $\mathcal{P}_\theta$ denote all data generating mechanisms $P$ satisfying the following conditions for known positive constants $M, B_L, B_U$.

(a) $\max\{|Y|, |A|\} \leq B_U$ a.s. $P$.
(b) $a \in H(\alpha, M)$, $b \in H(\beta, M)$, and $g \in H(\gamma, M)$.
(c) $0 < B_L < g(\mathbf{x}) < B_U$ for all $\mathbf{x} \in [0, 1]^d$.

Note that we do not put any assumptions on the function $c$. Indeed for $Y$ and $A$ binary random variables, the functions $a, b, g, c$ are variation independent. Following our discussion above, we will discuss adaptive estimation of $\phi(P) = \mathrm{E}_P(cov_P(Y, A|\mathbf{X})) = \mathrm{E}_P((Y - b(\mathbf{X}))(A - a(\mathbf{X})))$ over $P \in \mathcal{P}$. In particular, we summarize our results on upper and lower bounds on the adaptive minimax estimation of $\phi(P)$ in the following theorem.

THEOREM 3.1. *Assume* $(a) - (c)$ *and* $(\alpha, \beta, \gamma) \in \Theta$. *Then the following hold for positive* $C, C'$ *depending on* $M, B_L, B_U, \gamma_{\max}$.

(i) *(Upper Bound) There exists an estimator* $\hat{\phi}$, *depending only on* $M, B_L, B_U, \gamma_{\max}$ *such that*

$$\sup_{P \in \mathcal{P}_{(\alpha,\beta,\gamma)}} \mathrm{E}_P \left( \hat{\phi} - \phi(P) \right)^2 \leq C \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{4\alpha+4\beta}{2\alpha+2\beta+2d}}.$$

(ii) *(Lower Bound) Suppose* $\{A, Y\} \in \{0, 1\}^2$. *If one has*

$$\sup_{P \in \mathcal{P}_{(\alpha,\beta,\gamma)}} \mathrm{E}_P \left( \hat{\phi} - \phi(P) \right)^2 \leq C \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{4\alpha+4\beta}{2\alpha+2\beta+2d}},$$

*for an estimator* $\hat{\phi}$ *of* $\phi(P)$. *Then there exists a class of distributions* $\mathcal{P}_{(\alpha',\beta',\gamma')}$ *such that*

$$\sup_{P' \in \mathcal{P}_{(\alpha',\beta',\gamma')}} \mathrm{E}_{P'} \left( \hat{\phi} - \phi(P') \right)^2 \geq C' \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{4\alpha'+4\beta'}{2\alpha'+2\beta'+2d}}.$$

Theorem 3.1 describes the adaptive minimax estimation of the treatment effect functional in low regularity regime $(\frac{\alpha+\beta}{2} < d/4)$ i.e. when $\sqrt{n}$-rate estimation is not possible. By assuming more smoothness on the marginal of $g$ (i.e. lager value of $\gamma$) it is possible to include the case of $\frac{\alpha+\beta}{2} \geq d/4$ as well. In particular, if the set of $(\alpha, \beta)$ includes the case $\frac{\alpha+\beta}{2} > d/4$ as well as $\frac{\alpha+\beta}{2} < d/4$, one should be able to obtain adaptive and semi-parametrically efficient $\sqrt{n}$-consistent estimator of the treatment effect for $\frac{\alpha+\beta}{2} > d/4$. Similar to quadratic functionals (Giné and Nickl, 2008), the case of $\frac{\alpha+\beta}{2} = d/4$ will however incur an additional logarithmic penalty over usual $\sqrt{n}$-rate of

convergence. All these extensions can definitely be incorporated in the designing of the Lepski's method in Section 2. However, we do not pursue this in this paper and refer to Section 4 for more discussions on relevant issues. Finally, it is worth noting that if the set of $(\alpha, \beta)$ *only* includes the case $\frac{\alpha+\beta}{2} > d/4$, one can indeed obtain adaptive and even semiparametrically efficient estimation of the functionals studied here without effectively any assumption on $g$. The interested reader can find the details in Mukherjee, Newey and Robins (2017); Robins et al. (2016).

3.2. **Mean Response in Missing Data Models**. Suppose we have $n$ i.i.d observations on $O = (YA, A, \mathbf{X}) \sim P$, for a response variable $Y \in \mathbb{R}$ which is conditionally independent of the missingness indicator variable $A \in \{0, 1\}$ given covariate information $\mathbf{X}$. In literature, this assumption is typically known as the missing at random model (MAR) and under this assumption, our quantity of interest $\phi(P) = \mathrm{E}_P(Y)$ is identifiable as $\mathrm{EE}(Y|A = 1, X)$ from the observed data. This model is a canonical example of a study with missing response variable and to make this assumption reasonable, the covariates must contain the information on possible dependence between response and missingness. We referred the interested reader to Tsiatis (2007) for the history of statistical analysis of MAR and related models.

To the lay down the mathematical formalism for minimax adaptive estimation of $\phi(P)$ in this model, let $f$ be the marginal density of $\mathbf{X}$ (i.e. $\mathrm{E}_P(h(\mathbf{X})) = \int\limits_{[0,1]^d} h(\mathbf{x})f(\mathbf{x})d\mathbf{x}$ for all $P$-integrable function $h$), $a^{-1}(\mathbf{X}) := \mathrm{E}_P(A|\mathbf{X})$, and $b(\mathbf{X}) := \mathrm{E}_P(Y|A = 1, \mathbf{X}) = \mathrm{E}_P(Y|\mathbf{X})$, and $g(\mathbf{X}) = f(\mathbf{X})/a(\mathbf{X})$ (with the convention of the value $+\infty$ when dividing by 0). Although to be precise, we should put subscripts $P$ to $a, b, g$, we omit this since the context of their use is clear. Let $\Theta := \{\theta = (\alpha, \beta, \gamma) : \frac{\alpha+\beta}{2} < \frac{d}{4}, \gamma_{\max} \geq \gamma > \gamma_{\min} \geq 2(1 + \epsilon)\max\{\alpha, \beta\}\}$ for some fixed $\epsilon > 0$, and let $\mathcal{P}_\theta$ denote all data generating mechanisms $P$ satisfying the following conditions for known positive constants $M, B_L, B_U$.

(a) $|Y| \leq B_U$.
(b) $a \in H(\alpha, M)$, $b \in H(\beta, M)$, and $g \in H(\gamma, M)$.
(c) $B_L < g(\mathbf{x}), a(\mathbf{x}) < B_U$ for all $\mathbf{x} \in [0, 1]^d$.

We then have the following result.

THEOREM 3.2. *Assume $(a) - (c)$ and $(\alpha, \beta, \gamma) \in \Theta$. Then the following hold for positive $C, C'$ depending on $M, B_L, B_U, \gamma_{\max}$.*

*(i) (Upper Bound) There exists an estimator $\hat{\phi}$, depending only on $M, B_L, B_U, \gamma_{\max}$ such that*

$$\sup_{P \in \mathcal{P}_{(\alpha,\beta,\gamma)}} \mathrm{E}_P \left(\hat{\phi} - \phi(P)\right)^2 \leq C \left(\frac{\sqrt{\log n}}{n}\right)^{\frac{4\alpha+4\beta}{2\alpha+2\beta+2d}}.$$

*(ii) (Lower Bound) Suppose $\{A, Y\} \in \{0, 1\}^2$. If one has*

$$\sup_{P \in \mathcal{P}_{(\alpha,\beta,\gamma)}} \mathrm{E}_P \left(\hat{\phi} - \phi(P)\right)^2 \leq C \left(\frac{\sqrt{\log n}}{n}\right)^{\frac{4\alpha+4\beta}{2\alpha+2\beta+2d}},$$

*for an estimator $\hat{\phi}$ of $\phi(P)$. Then there exists a class of distributions $\mathcal{P}_{(\alpha',\beta',\gamma')}$ such that*

$$\sup_{P' \in \mathcal{P}_{(\alpha',\beta',\gamma')}} \mathrm{E}_{P'} \left(\hat{\phi} - \phi(P')\right)^2 \geq C' \left(\frac{\sqrt{\log n}}{n}\right)^{\frac{4\alpha'+4\beta'}{2\alpha'+2\beta'+2d}}.$$

Once again, Theorem 3.2 describes the adaptive minimax estimation of the average outcome in missing at random models in low regularity regime ($\frac{\alpha+\beta}{2} < d/4$) i.e. when $\sqrt{n}$-rate estimation is not possible. Extensions to include $\frac{\alpha+\beta}{2} \geq \frac{d}{4}$ is possible by similar Lepski type method with additional smoothness assumption on $g$. However, we do not pursue this in this paper and refer to Section 4 for more discussions on relevant issues.

3.3. **Quadratic and Variance Functionals in Regression Models.** Consider a observing data which are $n$ i.i.d copies of $O = (Y, \mathbf{X}) \sim P$ and the functional of interest is the expected value of the from of the regression of $Y$ on $\mathbf{X}$. Specifically suppose we want to estimate $\phi(P) = \mathrm{E}_P \left( \{E_P(Y|\mathbf{X})\}^2 \right)$. Assume that distribution of $\mathbf{X}$ has a density with respect to Lebesgue measure on $\mathbb{R}^d$ that has a compact support, which we assume to be $[0,1]^d$ for sake of simplicity. Let $g$ be the marginal density of $\mathbf{X}$, and $b(\mathbf{X}) := \mathrm{E}_P(Y|\mathbf{X})$. The class of distributions $\Theta := \{\mathcal{P}(\beta, \gamma) : \beta < \frac{d}{4}, \gamma_{\max} \geq \gamma > \gamma_{\min} \geq 2(1+\epsilon)\beta\}$ for some fixed $\epsilon > 0$, where by $\mathcal{P}(\beta, \gamma)$ we consider all data generating mechanisms $P$ satisfying the following conditions. for known positive constants $M, B_L, B_U$.

(a) $\max\{|Y|\} \leq B_U$.
(b) $b \in H(\beta, M)$, and $g \in H(\gamma, M)$.
(c) $0 < B_L < g(\mathbf{x}) < B_U$ for all $\mathbf{x} \in [0,1]^d$.

THEOREM 3.3. *Assume* $(a) - (c)$ *and* $(\beta, \gamma) \in \Theta$. *Then the following hold for positive* $C, C'$ *depending on* $M, B_L, B_U, \gamma_{\max}$.

(i) *(Upper Bound) There exists an estimator* $\hat{\phi}$, *depending only on* $M, B_L, B_U, \gamma_{\max}$ *such that*

$$\sup_{P \in \mathcal{P}(\beta, \gamma)} \mathrm{E}_P \left( \hat{\phi} - \phi(P) \right)^2 \leq C \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{8\beta}{4\beta+d}}.$$

(ii) *(Lower Bound) Suppose* $Y \in \{0,1\}^2$. *If one has*

$$\sup_{P \in \mathcal{P}(\beta, \gamma)} \mathrm{E}_P \left( \hat{\phi} - \phi(P) \right)^2 \leq C \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{8\beta}{4\beta+d}},$$

*for an estimator* $\hat{\phi}$ *of* $\phi(P)$. *Then there exists a class of distributions* $\mathcal{P}(\beta', \gamma')$ *such that*

$$\sup_{P' \in \mathcal{P}(\beta', \gamma')} \mathrm{E}_P \left( \hat{\phi} - \phi(P') \right)^2 \geq C' \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{8\beta'}{4\beta'+d}}.$$

REMARK 3.4. *Although Theorem 3.3 and the discussion before that is made in the context of estimating a particular quadratic functional in the context of a regression framework, it is worth noting that the result extends to estimating classical quadratic functionals in density models (Efromovich and Low, 1996; Giné and Nickl, 2008).*

One can also consider in the same set up, the estimation of functionals related to the conditional variance of $Y$ under such a regression model, which has been studied in detail by Brown and Levine (2007); Cai and Wang (2008); Fan and Yao (1998); Hall and Carroll (1989); Ruppert et al. (1997). Whereas, the minimax optimal and adaptive results in Brown and Levine (2007); Cai and Wang (2008) are in a equi-spaced fixed design setting, one can use an analogue of Theorem 3.3 to demonstrate a rate adaptive estimator and corresponding matching lower bound, with a mean-squared

error of the order of $\left(\frac{n}{\sqrt{\log n}}\right)^{-\frac{8\beta}{4\beta+d}}$ for estimating $\mathrm{E}_P(Var_P(Y|\mathbf{X}))$ adaptively over Hölder balls of regularity $\beta < \frac{d}{4}$. As noted by Robins et al. (2008), this rate is higher than the rate of estimating the conditional variance in mean-squared error for the equispaced design (Cai and Wang, 2008). In a similar vein, one can also obtain similar results for the estimation of conditional variance under the assumption if homoscedasticity i.e. $\sigma^2 := Var(Y|\mathbf{X} = \mathbf{x})$ for all $\mathbf{x} \in [0,1]^d$. In particular, one can once again obtain an estimator with mean-squared error of the order of $\left(\frac{n}{\sqrt{\log n}}\right)^{-\frac{8\beta}{4\beta+d}}$ for estimating $\sigma^2$ over any Hölder balls of regularity $\beta < \frac{d}{4}$. In particular, a candidate sequence of $\hat{\phi}_{n,j}$'s for this purpose was constructed in Robins et al. (2008), and in essence equals $\hat{\phi}_{n,j} = \frac{1}{n(n-1)}\sum_{i_1 \neq i_2}(Y_{i_1} - Y_{i_2})^2 K_{V_j}(\mathbf{X}_{i_1}, \mathbf{X}_{i_2})$— whose properties can be analyzed by a technique similar to the proof of Theorem 3.3. Its however worth noting that, even in this very easy to state classical problem, matching lower bounds, even in non-adaptive minimax sense is not available. Therefore, our results related to homoscedastic variance estimation, can be interpreted from a point of view of selecting a "best" candidate from a collection of estimators.

## 4. Discussion.
In this paper, we have extended the results for adaptive estimation of non-linear integral functionals from density estimation and Gaussian white noise models, to move towards an adaptive estimators non-linear functionals in more complex nonparametric models. Our results are provided and are most interesting in the low regularity regime, i.e. when it is not possible to achieve $\sqrt{n}$-efficient estimator in an asymptotically minimax sense. This typically happens when the "average smoothness" of the function classes in consideration is below $\frac{d}{4}$. The reason for focusing on the low regularity region is two fold. Firstly, this regime corresponds to situations where adaptation is to smoothness is not possible without paying a logarithmic penalty–making it an interesting case to study. Secondly, as noted in Robins et al. (2008), the appropriate non-adaptive minimax sequence of estimators of the functionals considered in this paper, which attain $\sqrt{n}$-efficiency rely either on high regularity of the marginal density of the covariates $\mathbf{X}$ in our examples or on correcting higher order bias by using U-statistics of degree 3 and higher. Indeed, under stringent assumptions on the smoothness of the density of $\mathbf{X}$, our results carry through to yield adaptive $\sqrt{n}$-efficient estimators. However, under more relaxed conditions on the density of $\mathbf{X}$, although we can in principle employ a similar framework of Lepski type idea as implemented by Theorem 2.1, the mathematical analysis of such a method requires sharp control of tails of U-statistics of degree 3 and higher. The structure of the higher order U-statistics considered in the estimators constructed in Robins et al. (2008) makes such an analysis delicate, and we plan to focus on these issues in a future paper. However, it is worth noting that with the additional knowledge of smoothness exceeding $\frac{d}{4}$ one can indeed obtain adaptive and even semiparametrically efficient estimation of the functionals studied here *without effectively any assumption on g*. The interested reader can find the details in Mukherjee, Newey and Robins (2017); Robins et al. (2016). Finally, the results of this paper can be extended to include the whole class of doubly robust functionals considered in Robins et al. (2008). However we only provide specific examples here to demonstrate the clear necessity to pay a poly-logarithmic penalty for adaptation in low regularity regimes.

## 5. Wavelets, Projections, and Hölder Spaces.
We work with certain Besov- Hölder type spaces which we define in terms of moduli of wavelet coefficients of continuous functions. For $d > 1$, consider expansions of functions $h \in L_2\left([0,1]^d\right)$ on an orthonormal basis of compactly supported bounded wavelets of the form

$$h(\mathbf{x}) = \sum_{k \in \mathbb{Z}^d} \langle h, \boldsymbol{\psi}_{0,k}^0 \rangle \boldsymbol{\psi}_{0,k}^0(\mathbf{x}) + \sum_{l=0}^{\infty} \sum_{k \in \mathbb{Z}^d} \sum_{v \in \{0,1\}^d - \{0\}^d} \langle h, \boldsymbol{\psi}_{l,k}^v \rangle \boldsymbol{\psi}_{l,k}^v(\mathbf{x}),$$

where the base functions $\boldsymbol{\psi}_{l,k}^v$ are orthogonal for different indices $(l,k,v)$ and are scaled and translated versions of the $2^d$ $S$-regular base functions $\boldsymbol{\psi}_{0,0}^v$ with $S > \beta$, i.e., $\boldsymbol{\psi}_{l,k}^v(x) = 2^{ld/2}\boldsymbol{\psi}_{0,0}^v(2^l\mathbf{x} - k) = \prod_{j=1}^d 2^{\frac{l}{2}}\psi_{0,0}^{v_j}\left(2^l x_j - k_j\right)$ for $k = (k_1,\ldots,k_d) \in \mathbb{Z}^d$ and $v = (v_1,\ldots,v_d) \in \{0,1\}^d$ with $\psi_{0,0}^0 = \phi$ and $\psi_{0,0}^1 = \psi$ being the scaling function and mother wavelet of regularity $S$ respectively as defined in one dimensional case. As our choices of wavelets, we will throughout use compactly supported scaling and wavelet functions of Cohen-Daubechies-Vial type with $S$ first null moments(Cohen, Daubechies and Vial, 1993). In view of the compact support of the wavelets, for each resolution level $l$ and index $v$, only $O(2^{ld})$ base elements $\psi_{l,k}^v$ are non-zero on $[0,1]$; let us denote the corresponding set of indices $k$ by $\mathcal{Z}_l$ obtaining the representation,

$$h(\mathbf{x}) = \sum_{k\in\mathcal{Z}_{J_0}} \langle h, \boldsymbol{\psi}_{J_0,k}^0\rangle \boldsymbol{\psi}_{J_0,k}^0(\mathbf{x}) + \sum_{l=J_0}^{\infty}\sum_{k\in\mathcal{Z}_l}\sum_{v\in\{0,1\}^d - \{0\}^d} \langle h, \boldsymbol{\psi}_{l,k}^v\rangle \boldsymbol{\psi}_{l,k}^v(\mathbf{x}),$$

(5.1)

where $J_0 = J_0(S) \geq 1$ is such that $2^{J_0} \geq S$ (Cohen, Daubechies and Vial, 1993; Giné and Nickl, 2015). Thereafter, let for any $h \in L_2[0,1]^d$, $\|\langle h, \boldsymbol{\psi}_{l',.}\rangle\|_2$ be the vector $L_2$ norm of the vector $\left(\langle h, \boldsymbol{\psi}_{l',k'}^v\rangle \colon k' \in \mathcal{Z}_{l'}, v \in \{0,1\}^d\right)$.

We will be working with projections onto subspaces defined by truncating expansions as above at certain resolution levels. For example letting

$$V_j := \text{span}\left\{\boldsymbol{\psi}_{l,k}^v, J_0 \leq l \leq j, k \in \mathcal{Z}_l, v \in \{0,1\}^d\right\}, j \geq J_0$$

(5.2)

one immediately has the following orthogonal projection kernel onto $V_j$ as

$$K_{V_j}(\mathbf{x}_1, \mathbf{x}_2) = \sum_{k\in\mathcal{Z}_{J_0}} \boldsymbol{\psi}_{J_0,k}^0(\mathbf{x}_1)\boldsymbol{\psi}_{J_0,k}^0(\mathbf{x}_2) + \sum_{l=J_0}^{j}\sum_{k\in\mathcal{Z}_l}\sum_{v\in\{0,1\}^d - \{0\}} \boldsymbol{\psi}_{l,k}^v(\mathbf{x}_1)\boldsymbol{\psi}_{l,k}^v(\mathbf{x}_2).$$

(5.3)

Owing to the MRA property of the wavelet basis, it is easy to see that $K_{V_j}$ has the equivalent representation as

$$K_{V_j}(\mathbf{x}_1, \mathbf{x}_2) = \sum_{k\in\mathcal{Z}_j}\sum_{v\in\{0,1\}^d} \psi_{jk}^v(x_1)\psi_{jk}^v(x_2).$$

(5.4)

Thereafter, using $S$-regular scaling and wavelet functions of Cohen-Daubechies-Vial type with $S > \beta$ let

$$H(\beta, M)$$
$$:= \left\{ \begin{array}{c} h \in C\left([0,1]^d\right) \\ : 2^{J_0(\beta+\frac{d}{2})}\|\langle h, \boldsymbol{\psi}_{J_0.}^0\rangle\|_\infty + \sup_{l\geq 0, k\in\mathbb{Z}^d, v\in\{0,1\}^d - \{0\}^d} 2^{l(\beta+\frac{d}{2})}|\langle h, \boldsymbol{\psi}_{l,k}^v\rangle| \leq M \end{array} \right\},$$

(5.5)

with $C\left([0,1]^d\right)$ being the set of all continuous bounded functions on $[0,1]^d$. It is standard result in the theory of wavelets that $H(\beta, M)$ is related to classical Hölder-Zygmund spaces with equivalent norms (see (Giné and Nickl, 2015, Chapter 4) for details). For $0 < \beta < 1$ for example, $H(\beta, M)$

consists of all functions in $C\left([0,1]^d\right)$ such that $\|f\|_\infty + \sup\limits_{\mathbf{x}_1, \mathbf{x}_2 \in [0,1]^d} \frac{|f(\mathbf{x}_1) - f(\mathbf{x}_2)|}{\|\mathbf{x}_1 - \mathbf{x}_2\|^\beta} \leq C(M)$. For non-integer $\beta > 1$, $H(\beta, M)$ consists of all functions in $C\left([0,1]^d\right)$ such that $f^{(\lfloor\beta\rfloor)} \in C\left([0,1]^d\right)$ for any partial $f^{(\lfloor\beta\rfloor)}$ of order $\lfloor\beta\rfloor$ of $f$ and $\|f\|_\infty + \sup\limits_{\mathbf{x}_1, \mathbf{x}_2 \in [0,1]^d} \frac{|f^{(\lfloor\beta\rfloor)}(\mathbf{x}_1) - f^{(\lfloor\beta\rfloor)}(\mathbf{x}_2)|}{\|\mathbf{x}_1 - \mathbf{x}_2\|^{\beta - \lfloor\beta\rfloor}} \leq C(M)$. Therefore, the functions in $H(\beta, M)$ are automatically uniformly bounded by a number depending on the radius $M$.

## 6. Proof of Main Theorems.

*Proof of Theorem 2.1.*

PROOF. In this proof we repeatedly use the fact that for any fixed $m \in \mathbb{N}$ and $a_1, \ldots, a_m$ real numbers, one has by Hölder's Inequality $|a_1 + \ldots, a_m|^p \leq C(m, p)\left(|a_1|^p + \ldots + |a_m|^p\right)$ for $p > 1$. Suppose $\beta \in (\beta_{l+1}, \beta_l]$ for some $l = 0, \ldots, N-2$. Indeed, this $l$ depends on $\beta$ and $n$, and therefore, to be very precise, we should write it as $l(\beta, n)$. However, for the sake of brevity we omit such notation. We immediately have the following lemma.

LEMMA 6.1. *For $l' \geq l+2$ and $C_{\mathrm{opt}}$ large enough depending on $C_1, C_2, B, \psi_{0,0}^0, \psi_{0,0}^1$,*

$$\sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta) = \beta, f_2(\theta) > \frac{4\beta}{4\beta + d}}} \mathbb{P}_P\left(\hat{l} = l'\right) \leq \frac{C}{n},$$

*where $C > 0$ is an universal constant.*

PROOF.

$$\sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta) = \beta, f_2(\theta) > \frac{4\beta}{4\beta + d}}} \mathbb{P}_P\left(\hat{l} \geq l+2\right)$$

$$\leq \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta) = \beta, f_2(\theta) > \frac{4\beta}{4\beta + d}}} \mathbb{P}_P\left(\begin{array}{c} \exists l' > l+1: \\ \left(\hat{\phi}_{n, j(k_*(j_{l+1}))} - \hat{\phi}_{n, j(k_*(j_{l'}))}\right)^2 > C_{\mathrm{opt}}^2 \log nR(k^*(j_{l'})) \end{array}\right)$$

$$\leq \sum_{l'=l+1}^{N} \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta) = \beta, f_2(\theta) > \frac{4\beta}{4\beta + d}}} \mathbb{P}_P\left(|\hat{\phi}_{n, j(k_*(j_{l+1}))} - \hat{\phi}_{n, j(k_*(j_{l'}))}| > C_{\mathrm{opt}}\sqrt{\log nR(k^*(j_{l'}))}\right)$$

For any fixed $l+2 \leq l' \leq N_1$, using that $R(k_*(j_{l+1})) \leq R(k_*(j_{l'}))$ and

$$\sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta) = \beta, f_2(\theta) > \frac{4\beta}{4\beta + d}}} |\mathrm{E}_P\left(\hat{\phi}_{n, j(k_*(j_{l+1}))}\right) - \mathrm{E}_P\left(\hat{\phi}_{n, j(k_*(j_{l'}))}\right)|$$

$$\leq Ck_*(j_{l+1})^{-2\beta_{l+1}/d} + \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta) = \beta, f_2(\theta) > \frac{4\beta}{4\beta + d}}} n^{-f_2(\theta)} \leq 2C\sqrt{\log nR(k^*(j_{l+1}))}$$

we have that

$$\sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta) = \beta, f_2(\theta) > \frac{4\beta}{4\beta + d}}} \mathbb{P}_P\left(|\hat{\phi}_{n, j(k_*(j_{l+1}))} - \hat{\phi}_{n, j(k_*(j_{l'}))}| > C_{\mathrm{opt}}\sqrt{\log nR(k^*(j_{l'}))}\right)$$

$$
\leq \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathbb{P}_P \left( |\hat{\phi}_{n,j(k_*(j_{l+1}))} - \mathrm{E}_P \left( \hat{\phi}_{n,j(k_*(j_{l+1}))} \right) | > \frac{C_{\mathrm{opt}}}{2} \sqrt{\log n R(k^*(j_{l'}))} \right)
$$

$$
+ \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathbb{P}_P \left( \begin{array}{c} |\hat{\phi}_{n,j(k_*(j_{l'}))} - \mathrm{E}_P \left( \hat{\phi}_{n,j(k_*(j_{l'}))} \right) | > \frac{C_{\mathrm{opt}}}{2} \sqrt{\log n R(k^*(j_{l'}))} \\ -|\mathrm{E}_P \left( \hat{\phi}_{n,j(k_*(j_{l+1}))} \right) - \mathrm{E}_P \left( \hat{\phi}_{n,j(k_*(j_{l'}))} \right) | \end{array} \right)
$$

$$
\leq \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathbb{P}_P \left( |\hat{\phi}_{n,j(k_*(j_{l+1}))} - \mathrm{E}_P \left( \hat{\phi}_{n,j(k_*(j_{l+1}))} \right) | > \frac{C_{\mathrm{opt}}}{2} \sqrt{\log n R(k^*(j_{l+1}))} \right)
$$

$$
+ \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathbb{P}_P \left( \begin{array}{c} |\hat{\phi}_{n,j(k_*(j_{l'}))} - \mathrm{E}_P \left( \hat{\phi}_{n,j(k_*(j_{l'}))} \right) | > \frac{C_{\mathrm{opt}}}{2} \sqrt{\log n R(k^*(j_{l'}))} \\ -|\mathrm{E}_P \left( \hat{\phi}_{n,j(k_*(j_{l+1}))} \right) - \mathrm{E}_P \left( \hat{\phi}_{n,j(k_*(j_{l'}))} \right) | \end{array} \right)
$$

Now

$$
\sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} |\mathrm{E}_P \left( \hat{\phi}_{n,j(k_*(j_{l+1}))} \right) - \mathrm{E}_P \left( \hat{\phi}_{n,j(k_*(j_{l'}))} \right) |
$$

$$
\leq \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} 2C_1 \left( 2^{-2j(k_*(j_{l+1}))f_1(\theta)} + n^{-f_2(\theta)} \right)
$$

$$
\leq 2C_1 n^{-\frac{4\beta}{4\beta+d}} + 2C_1 \left( \frac{k_*(j_{l+1})}{2^d} \right)^{-2\beta/d} \leq 2C_1 n^{-\frac{4\beta}{4\beta+d}} + 2C_1 2^{d\beta} \left( \frac{n}{\sqrt{\log n}} \right)^{-\frac{4\beta}{4\beta_{l+1}+d}}
$$

$$
\leq 2C_1 n^{-\frac{4\beta}{4\beta+d}} + 2C_1 2^{d^2/4} \left( \frac{n}{\sqrt{\log n}} \right)^{-\frac{4\beta}{4\beta_{l+1}+d}}.
$$

The last quantity in the above display is smaller than $\frac{C_{\mathrm{opt}}}{4} \sqrt{\log n R(k^*(j_{l'}))}$ for $C_{\mathrm{opt}}$ chosen large enough (depending on $C_1$ and $d$). This implies that for $C_{\mathrm{opt}}$ properly chosen based on the given parameters,

$$
\sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathbb{P}_P \left( |\hat{\phi}_{n,j(k_*(j_{l+1}))} - \hat{\phi}_{n,j(k_*(j_{l'}))}| > C_{\mathrm{opt}} \sqrt{\log n R(k^*(j_{l'}))} \right)
$$

$$
\leq \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathbb{P}_P \left( |\hat{\phi}_{n,j(k_*(j_{l+1}))} - \mathrm{E}_P \left( \hat{\phi}_{n,j(k_*(j_{l+1}))} \right) | > \frac{C_{\mathrm{opt}}}{2} \sqrt{\log n R(k^*(j_{l+1}))} \right)
$$

$$
+ \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathbb{P}_P \left( |\hat{\phi}_{n,j(k_*(j_{l'}))} - \mathrm{E}_P \left( \hat{\phi}_{n,j(k_*(j_{l'}))} \right) | > \frac{C_{\mathrm{opt}}}{4} \sqrt{\log n R(k^*(j_{l'}))} \right)
$$

$$
\leq \frac{C}{n},
$$

for an universal constant $C$. The last inequality can now be obtained by standard Hoeffding's decomposition and subsequent application of Lemma B.2 and B.5 to the second and first order degenerate parts respectively. The control thereafter is standard by our choices of $n \leq 2^{jd} \leq n^2$. For similar calculations we refer to (Mukherjee and Sen, 2016). ∎

Returning to the proof of Theorem 2.1 we have

$$\sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathrm{E}_P \left( \hat{\phi}_{n,j\left(k^*(j_{\hat{l}})\right)} - \phi(P) \right)^2$$

$$\leq \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathrm{E}_P \left( \left( \hat{\phi}_{n,j\left(k^*(j_{\hat{l}})\right)} - \phi(P) \right)^2 \mathcal{I}\left( \hat{l} \leq l+1 \right) \right)$$

$$+ \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathrm{E}_P \left( \left( \hat{\phi}_{n,j\left(k^*(j_{\hat{l}})\right)} - \phi(P) \right)^2 \mathcal{I}\left( \hat{l} \geq l+2 \right) \right)$$

$$= T_1 + T_2$$

*Control of $T_1$.* Using the definition of $\hat{l}$ we have the following string of inequalities.

$$T_1$$

$$= \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathrm{E}_P \left[ \mathcal{I}\left( \hat{l} \leq l+1 \right) \left( \hat{\phi}_{n,j\left(k^*(j_{\hat{l}})\right)} - \phi(P) \right)^2 \right]$$

$$\leq 2 \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathrm{E}_P \left[ \mathcal{I}\left( \hat{l} \leq l+1 \right) \left\{ \begin{array}{c} \left( \hat{\phi}_{n,j\left(k^*(j_{\hat{l}})\right)} - \hat{\phi}_{n,j(k^*(j_{l+1}))} \right)^2 \\ + \left( \hat{\phi}_{n,j(k^*(j_{l+1}))} - \phi(P) \right)^2 \end{array} \right\} \right]$$

$$\leq 2 \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} C_{\mathrm{opt}}^2 \log n R\left( k^*(j_{l+1}) \right) + 2 \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathrm{E}_P \left[ \left( \hat{\phi}_{n,j(k^*(j_{l+1}))} - \phi(P) \right)^2 \right]$$

$$\leq C_{\mathrm{opt}}^2 \log n R\left( k^*(j_{l+1}) \right) + C \frac{2^{j(k_*(j_{l+1}))d}}{n^2} + C 2^{-4j(k_*(j_{l+1}))d\frac{\beta}{d}} + \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} 2Cn^{-2f_2(\theta)}$$

$$\leq (4C + C_{\mathrm{opt}}^2) \left( \frac{n^2}{\log n} \right)^{-\frac{4\beta_{l+1}/d}{1+4\beta_{l+1}/d}} \leq e^{c'} \left( \frac{n^2}{\log n} \right)^{-\frac{4\beta/d}{1+4\beta/d}}.$$

Above we have used the property of $r(\beta)$, definition of $j(k^*(l))$ to conclude that $2^{j(k^*(l))d} \leq k^*(l)$, and also the definition and properties of $\beta_l$ together with the fact that $\beta_{l+1} \leq \beta < \beta_l$ implies $\beta/d - \beta_{l+1}/d \leq \frac{c}{\log n}$ for some fixed constant $c$.

*Control of $T_2$.*

$$T_2$$

$$\leq \sum_{l'=l+2}^{N} \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathrm{E}_P \left[ \mathcal{I}(\hat{l} = l') \left( \hat{\phi}_{n,j(k^*(j_{l'}))} - \phi(P) \right)^2 \right]$$

$$\leq \sum_{l'=l+2}^{N} \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathbb{P}_P^{\frac{1}{p}} \left( \hat{l} = l' \right) \sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\beta, f_2(\theta) > \frac{4\beta}{4\beta+d}}} \mathrm{E}_P^{\frac{1}{q}} \left[ \left( \hat{\phi}_{n,j\left(k^*(j_{l'})\right)} - \phi(P) \right)^{2q} \right]$$

$$\leq \sum_{l'=j+2}^{N} 4 \left( \frac{C}{n} \right)^{\frac{1}{p}} \left[ \frac{2^{j\left(k^*(j_{l'})\right)d}}{n^2} + 2^{-4j\left(k^*(j_{l'})\right)\beta/d} + n^{-\frac{8\beta}{4\beta+d}} \right] \lesssim \log n \left( \frac{C}{n} \right)^{\frac{1}{p}},$$

where the second last inequality above follows from Lemma 6.1, Lemma B.2, and Lemma B.5, and the last inequality follows from the choice of $N \lesssim \log n$. Therefore for $p$ sufficiently close to 1, we have desired control over $T_2$.

∎

*Proof of Theorem 2.2.*

PROOF. We closely follow the proof strategy of Robins et al. (2009). In particular, consider $\overline{P}_n := \int P_\lambda^n d\pi(\lambda)$ and $\overline{Q}_n := \int Q_\lambda^n d\pi(\lambda)$ to be distributions of $O_1, \ldots, O_n$ and define probability measures on $\chi_j$ defined by $P_{j,\lambda_j} = \frac{\mathcal{I}_{\chi_j} dP_\lambda}{p_j}$ and $Q_{j,\lambda_j} = \frac{\mathcal{I}_{\chi_j} dQ_\lambda}{p_j}$. Therefore, if $p_\lambda$ and $q_\lambda$ are densities of $P_\lambda$ and $Q_\lambda$ respectively with respect to some common dominating measure, then one has

$$\chi^2 \left( \int P_\lambda^n d\pi(\lambda), \int Q_\lambda^n d\pi(\lambda) \right)$$

$$= \mathrm{E}_{\overline{P}_n} \left( \frac{\int \prod\limits_{j=1}^{k} \prod\limits_{i} \mathcal{I}(O_i \in \chi_j) q_\lambda(O_i) d\pi_j(\lambda_j)}{\int \prod\limits_{j=1}^{k} \prod\limits_{i} \mathcal{I}(O_i \in \chi_j) p_\lambda(O_i) d\pi_j(\lambda_j)} \right)^2 - 1$$

$$= \mathrm{E}_{\overline{P}_n} \left( \prod_{j=1}^{k} \frac{\int \prod\limits_{i:O_i \in \chi_j} q_{j,\lambda_j}(O_i) d\pi_j(\lambda_j)}{\int \prod\limits_{i:O_i \in \chi_j} p_{j,\lambda_j}(O_i) d\pi_j(\lambda_j)} \right)^2 - 1. \tag{6.1}$$

Define variables $I_1, \ldots, I_n$ such that $I_i = j$ if $O_I \in \chi_j$ for $i = 1, \ldots, n$, $j = 1, \ldots, k$ and let $N_j = |\{i: I_i = j\}|$. Now note that the measure $\overline{P}_n$ arises as the distribution $O_1, \ldots, O_n$ if this vector is generated in two steps as follows. First one chooses $\lambda \sim \pi$ and then given $\lambda$ $O_1, \ldots, O_n$ are generated independent $p_\lambda$. This implies that given $\lambda$, $(N_1, \ldots, N_k)$ is distributed as Multinomial$(P_\lambda(\chi_1), \ldots, P_\lambda(\chi_k))$ = Multinomial$(p_1, \ldots, p_k) \Rightarrow (N_1, \ldots, N_k) \perp \lambda \Rightarrow$ under $\overline{P}_n$, $(N_1, \ldots, N_k) \sim$ Multinomial$(p_1, \ldots, p_k)$ unconditionally. Similarly, given $\lambda$, $I_1, \ldots, I_n$ are independent, and the event $I_i = j$ has probability $p_j$ which is again free of $\lambda$. This in turn implies that $(I_1, \ldots, I_n) \perp \lambda$ under $\overline{P}_n$. The conditional distribution of $O_1, \ldots, O_n$ given $\lambda$ and $(I_1, \ldots, I_n)$ can also be described as follows. For each partitioning set $\chi_j$ generate $N_j$ variables independently from $P_\lambda$ restricted and renormalized to $\chi_j$, i.e. from the measure $P_{j,\lambda_j}$. Now we can do so independently across the partitioning sets and attach correct labels $\{1, \ldots, n\}$ which are consistent with $(I_1, \ldots, I_n)$. The conditional distribution of $O_1, \ldots, O_n$ under $\overline{P}_n$ given $I_1, \ldots, I_n$ is the mixture of this distribution relative to the conditional distribution of $\lambda$ given $I_1, \ldots, I_n$, which by the virtue of independence of $I_1, \ldots, I_n$ and $\lambda$ under $\overline{P}_n$, was seen to be the unconditional distribution, $\pi$. Thus we can obtain a sample from the conditional distribution under $\overline{P}_n$ of $O_1, \ldots, O_n$ given $I_1, \ldots, I_n$ by generating for each partitioning set $\chi_j$ a set of $N_j$ variables from the measure $\int P_{j,\lambda_j} d\pi(\lambda_j)$, independently across the partitioning sets, and next attaching labels consistent with $I_1, \ldots, I_n$. The above discussion allows us to write the right hand side of (6.1) as

$$\mathrm{E}_{\overline{P}_n} \mathrm{E}_{\overline{P}_n} \left[ \prod_{j=1}^{k} \left( \frac{\int \prod\limits_{i:I_i=j} q_{j,\lambda_j}(O_i) d\pi_j(\lambda_j)}{\int \prod\limits_{I_i=j} p_{j,\lambda_j}(O_i) d\pi_j(\lambda_j)} \right)^2 | I_1, \ldots, I_n \right] - 1$$

$$= \mathrm{E}_{\overline{P}_n} \prod_{j=1}^{k} \mathrm{E}_{\overline{P}_n} \left[ \left( \frac{\int \prod\limits_{i:I_i=j} q_{j,\lambda_j}(O_i) d\pi_j(\lambda_j)}{\int \prod\limits_{I_i=j} p_{j,\lambda_j}(O_i) d\pi_j(\lambda_j)} \right)^2 | I_1, \ldots, I_n \right] - 1$$

$$= \mathrm{E}\left[\prod_{j=1}^{k}\mathrm{E}_{\int P_{j,\lambda_j}^{N_j}d\pi_j(\lambda_j)}\left(\frac{\int Q_{j,\lambda_j}^{N_j}d\pi_j(\lambda_j)}{\int P_{j,\lambda_j}^{N_j}d\pi_j(\lambda_j)}\right)^2\right]-1$$

Arguing similar to Lemma 5.2 of Robins et al. (2009), we have with $\widetilde{c}=\max_j\sup_\lambda\int_{\chi_j}\frac{p^2}{p_\lambda}\frac{d\mu}{p_j}$ and $c=\max\{\widetilde{c},1\}$ that

$$\mathrm{E}\left[\prod_{j=1}^{k}\mathrm{E}_{\int P_{j,\lambda_j}^{N_j}d\pi_j(\lambda_j)}\left(\frac{\int Q_{j,\lambda_j}^{N_j}d\pi_j(\lambda_j)}{\int P_{j,\lambda_j}^{N_j}d\pi_j(\lambda_j)}\right)^2\right]-1$$

$$\leq \mathrm{E}\left[\prod_{j=1}^{k}\left[1+2\left(\sum_{r=2}^{N_j}\binom{N_j}{r}b^r+2N_j^2\sum_{r=1}^{N_j-1}\binom{N_j-1}{r}a^rb+2N_j^2\widetilde{c}^{N_j-1}\delta\right)\right]\right]-1$$

$$= \mathrm{E}\left[\prod_{j=1}^{k}\left[1+2\left(((1+b)^{N_j}-1-N_jb)+N_j^2\left((1+a)^{N_j-1}-1\right)b+2N_j^2\widetilde{c}^{N_j-1}\delta\right)\right]\right]-1$$

$$\leq \mathrm{E}\left[\prod_{j=1}^{k}\left[1+2\left(((1+b)^{N_j}-1-N_jb)+N_j^2\left((1+a)^{N_j-1}-1\right)b+2N_j^2c^{N_j-1}\delta\right)\right]\right]-1$$

To bound the last quantity in the above display we use Shao (2000) to note that since $N_j, j = 1,\ldots,k$ is a multinomial random vector, for any increasing function $f$ on the range of $N_j$'s, $\mathrm{E}\left(\prod_{j=1}^{k}f(N_j)\right)=\mathrm{E}\left(\exp\left(\sum_{j=1}^{k}N_j\right)\right)\leq\exp\left(\sum_{j=1}^{k}\mathrm{E}f(N_j)\right)=\prod_{j=1}^{k}\mathbb{E}_f(N_j)$. In this context, noting that

$$1+2\left(((1+b)^{N_j}-1-N_jb)+N_j^2\left((1+a)^{N_j-1}-1\right)b+2N_j^2c^{N_j-1}\delta\right)$$

is an increasing function in each coordinate of $N_j$ (on the range of $N_j$) we therefore have by Shao (2000) that the last display is bounded by

$$\prod_{j=1}^{k}\mathrm{E}\left[\left[1+2\left(((1+b)^{N_j}-1-N_jb)+N_j^2\left((1+a)^{N_j-1}-1\right)b+2N_j^2c^{N_j-1}\delta\right)\right]\right]-1$$

$$= \prod_{j=1}^{k}\left[1+2\left\{\begin{array}{c}(1+bp_j)^n-1-nbp_j\\+np_j(1+ap_j)^{n-2}(1+nap_j+np_j-p_j)b-np_j(1-p_j)b-np_j^2b\\+2\delta np_j(cp_j+1-p_j)^{n-2}(cnp_j+1-p_j)\end{array}\right\}\right]-1$$

$$\leq \prod_{j=1}^{k}\left[1+C\left((np_jb)^2+(np_j)^2ab+np_j\delta\right)\right]-1$$

$$\leq \prod_{j=1}^{k}\left[1+Cn(\max_j p_j)\left((np_j)b^2+np_jab\right)+np_j\delta\right]-1$$

$$\leq e^{\sum_{j=1}^{k}\left(Cn(\max_j p_j)\left((np_j)b^2+np_jab\right)+np_j\delta\right)}-1=e^{Cn^2(\max_j p_j)(b^2+ab)+Cn\delta}-1$$

where we have used the fact that $n(\max_j p_j)(1\vee a\vee b\vee\widetilde{c})\leq A$ for a positive constant $A$ along with $\sum_{j=1}^{k}p_j=1$.

∎

6.1. *Proof of Theorem 2.3.*

PROOF. Let

$$2^{j_{\min}d} = \left\lfloor \left( \frac{n}{\log n} \right)^{\frac{1}{2\beta_{\max}/d+1}} \right\rfloor, \quad 2^{j_{\max}d} = \left\lfloor \left( \frac{n}{\log n} \right)^{\frac{1}{2\beta_{\min}/d+1}} \right\rfloor,$$

$$2^{l_{\min}d} = \left\lfloor \left( \frac{n}{\log n} \right)^{\frac{1}{2\gamma_{\max}/d+1}} \right\rfloor, \quad 2^{l_{\max}d} = \left\lfloor \left( \frac{n}{\log n} \right)^{\frac{1}{2\gamma_{\min}/d+1}} \right\rfloor.$$

Without loss of generality assume that we have data $\{\mathbf{x}_i, y_i\}_{i=1}^{2n}$. We split it into two equal parts and use the second part to construct the estimator $\hat{g}$ of the design density $g$ and us the resulting $\hat{g}$ to construct the adaptive estimate of the regression function from the first half of the sample. Throughout the proof, $E_{P,i}[\cdot]$ will denote the expectation with respect to the $i^{th}$ half of the sample, with the other half held fixed, under the distribution $P$. Throughout we choose the regularity of our wavelet bases to be larger than $\gamma_{\max}$ for the desired approximation and moment properties to hold. As a result our constants depend on $\gamma_{\max}$.

Define $\mathcal{T}_1 = [j_{\min}, j_{\max}] \cap \mathbb{N}$ and $\mathcal{T}_2 = [l_{\min}, l_{\max}] \cap \mathbb{N}$. For $l \in \mathcal{T}_2$, let $\hat{g}_l(\mathbf{x}) = \frac{1}{n} \sum_{i=n+1}^{2n} K_{V_l}(\mathbf{X}_i, \mathbf{x})$. Now, let

$$\hat{l} = \min \left\{ j \in \mathcal{T}_2 \colon \|\hat{g}_j - \hat{g}_l\|_\infty \le C^* \sqrt{\frac{2^{ld} ld}{n}}, \ \forall l \in \mathcal{T}_2 \text{ s.t. } l \ge j \right\}.$$

where $C^*$ is a constant (depending on $\gamma_{\max}, B_U$) that can be determined from the proof hereafter. Thereafter, consider the estimator $\widetilde{g} := \hat{g}_{\hat{l}}$.

Fix a $P := (f, g) \in \mathcal{P}(\beta, \gamma)$. To analyze the estimator $\widetilde{g}$, we begin with standard bias variance type analysis for the candidate estimators $\hat{g}_l$. First note that for any $\mathbf{x} \in [0,1]^d$, using standard facts about compactly supported wavelet basis having regularity larger than $\gamma_{\max}$ (Härdle et al., 1998), one has for a constant $C_1$ depending only on $q$ and the wavelet basis used,

$$|\mathrm{E}_P(\hat{g}_l(\mathbf{x})) - g(\mathbf{x})| = |\Pi(g|V_l)(\mathbf{x}) - g(\mathbf{x})| \le C_1 M 2^{-ld\frac{\gamma}{d}}. \tag{6.2}$$

Above we have used the fact that

$$\sup_{h \in H(\gamma, M)} \|h - \Pi(h|V_l)\|_\infty \le C_1 M 2^{-l\gamma}. \tag{6.3}$$

Also, by standard arguments about compactly supported wavelet basis having regularity larger than $\gamma_{\max}$ (Giné and Nickl, 2015), one has for a constant $C_2 := C(B_U, \boldsymbol{\psi}_{0,0}^0, \boldsymbol{\psi}_{0,0}^1, \gamma_{\max})$

$$\mathrm{E}_P\left( \|\hat{g}_l(\mathbf{x}) - \mathrm{E}_P(\hat{g}_l(\mathbf{x}))\|_\infty \right) \le C_2 \sqrt{\frac{2^{ld} ld}{n}}. \tag{6.4}$$

Therefore, by (6.2), (6.4), and triangle inequality,

$$\mathrm{E}_{P,2} \|\hat{g}_l - g\|_\infty \le C_1 M 2^{-ld\frac{\gamma}{d}} + C_2 \sqrt{\frac{2^{ld} ld}{n}}.$$

Define,

$$l^* := \min \left\{ l \in \mathcal{T}_2 \colon C_1 M 2^{-ld\frac{\gamma}{d}} \le C_2 \sqrt{\frac{2^{ld} ld}{n}} \right\}.$$

The definition of $l^*$ implies that for $n$ sufficiently large,

$$2^d \left( \frac{C_1}{C_2} M \right)^{\frac{2d}{2\gamma+d}} \left( \frac{n}{\log n} \right)^{\frac{d}{2\gamma+d}} \leq 2^{l^* d} \leq 2^{d+1} \left( \frac{C_1}{C_2} M \right)^{\frac{2d}{2\gamma+d}} \left( \frac{n}{\log n} \right)^{\frac{d}{2\gamma+d}} . \tag{6.5}$$

The error analysis of $\widetilde{g}$ can now be carried out as follows.

$$\mathrm{E}_{P,2}\|\widetilde{g} - g\|_\infty = \mathrm{E}_{P,2}\|\widetilde{g} - g\|_\infty \mathcal{I}\left( \hat{l} \leq l^* \right) + \mathrm{E}_{P,2}\|\widetilde{g} - g\|_\infty \mathcal{I}\left( \hat{l} > l^* \right)$$
$$:= I + II. \tag{6.6}$$

We first control term $I$ as follows.

$$I = \mathrm{E}_{P,2}\|\widetilde{g} - g\|_\infty \mathcal{I}\left( \hat{l} \leq l^* \right)$$
$$\leq \mathrm{E}_{P,2}\|\hat{g}_{\hat{l}} - \hat{g}_{l^*}\|_\infty \mathcal{I}\left( \hat{l} \leq l^* \right) + \mathrm{E}_{P,2}\|\hat{g}_{l^*} - g\|_\infty \mathcal{I}\left( \hat{l} \leq l^* \right)$$
$$\leq C^* \sqrt{\frac{2^{l^* d} l^* d}{n}} + C_1 M 2^{-l^* d \frac{\gamma}{d}} + C_2 \sqrt{\frac{2^{l^* d} l^* d}{n}}$$
$$\leq (C^* + 2C_2) \sqrt{\frac{2^{l^* d} l^* d}{n}} \leq 2^{d+1} \left( \frac{C_1}{C_2} M \right)^{\frac{2d}{2\gamma+d}} \left( \frac{n}{\log n} \right)^{-\frac{\gamma}{2\gamma+d}} . \tag{6.7}$$

The control of term $II$ is easier if one has suitable bounds on $\|\hat{g}_l - g\|_\infty$. To this end note that, for any fixed $\mathbf{x} \in [0,1]^d$, there exists a constant $C_3 := C(\boldsymbol{\psi}_{0,0}^0, \boldsymbol{\psi}_{0,0}^1, \gamma_{\max})$

$$|\hat{g}_l(\mathbf{x})| \leq \frac{1}{n} \sum_{i=1}^n \sum_{k \in \mathcal{Z}_l} \sum_{v \in \{0,1\}^d} |\psi_{l,k}^v(\mathbf{X}_i)| |\psi_{l,k}^v(\mathbf{x})| \leq C_3 2^{ld}.$$

This along with the fact that $\|g\|_\infty \leq B_U$, implies that for $n$ sufficiently large,

$$\|\hat{g}_l - g\|_\infty \leq C_3 2^{ld} + B_U \leq 2C_3 2^{ld}.$$

In the above display the last inequality follows since $l \geq l_{\min} \geq \left( \frac{n}{\log n} \right)^{\frac{1}{2\gamma_{\max}/d+1}}$. Therefore,

$$II \leq C_3 \sum_{l=\hat{l}}^{l_{\max}} 2^{ld} \mathbb{P}\left( \hat{l} = l \right). \tag{6.8}$$

We now complete the control over II by suitably bounding $\mathbb{P}\left( \hat{l} = l \right)$. To this end, note that for any $l > l^*$,

$$\mathbb{P}_{P,2}\left( \hat{l} = l \right)$$
$$\leq \sum_{l>l^*} \mathbb{P}_{P,2}\left( \|\hat{g}_l - \hat{g}_{l^*}\|_\infty > C^* \sqrt{\frac{2^{ld} l d}{n}} \right)$$
$$\leq \sum_{l>l^*} \left\{ \begin{array}{c} \mathbb{P}_{P,2}\left( \|\hat{g}_{l^*} - \mathrm{E}\left( \hat{g}_{l^*} \right)\|_\infty > \frac{C^*}{2}\sqrt{\frac{2^{ld} l d}{n}} - \|\mathrm{E}_{P,2}\left( \hat{g}_{l^*} \right) - \mathrm{E}_{P,2}\left( \hat{g}_l \right)\|_\infty \right) \\ +\mathbb{P}_{P,2}\left( \|\hat{g}_l - \mathrm{E}_{P,2}\left( \hat{g}_l \right)\|_\infty > \frac{C^*}{2}\sqrt{\frac{2^{ld} l d}{n}} \right) \end{array} \right\}$$

$$\leq \sum_{l>l^*} \left\{ \begin{array}{c} \mathbb{P}_{P,2}\left( \|\hat{g}_{l^*} - \mathrm{E}\left(\hat{g}_{l^*}\right)\|_\infty > \frac{C^*}{2}\sqrt{\frac{2^{ld}ld}{n}} - \|\Pi\left(g|V_{l^*}\right) - \Pi\left(g|V_l\right)\|_\infty \right) \\ + \mathbb{P}\left( \|\hat{g}_l - \mathrm{E}_{P,2}\left(\hat{g}_l\right)\|_\infty > \frac{C^*}{2}\sqrt{\frac{2^{ld}ld}{n}} \right) \end{array} \right\}$$

$$\leq \sum_{l>l^*} \left\{ \begin{array}{c} \mathbb{P}_{P,2}\left( \|\hat{g}_{l^*} - \mathrm{E}\left(\hat{g}_{l^*}\right)\|_\infty > \frac{C^*}{2}\sqrt{\frac{2^{ld}ld}{n}} - 2C_2\sqrt{\frac{2^{l^*d}l^*d}{n}} \right) \\ + \mathbb{P}\left( \|\hat{g}_l - \mathrm{E}_{P,2}\left(\hat{g}_l\right)\|_\infty > \frac{C^*}{2}\sqrt{\frac{2^{ld}ld}{n}} \right) \end{array} \right\}$$

$$\leq \sum_{l>l^*} \left\{ \begin{array}{c} \mathbb{P}_{P,2}\left( \|\hat{g}_{l^*} - \mathrm{E}\left(\hat{g}_{l^*}\right)\|_\infty > (\frac{C^*}{2} - 2C_2)\sqrt{\frac{2^{ld}ld}{n}} \right) \\ + \mathbb{P}\left( \|\hat{g}_l - \mathrm{E}_{P,2}\left(\hat{g}_l\right)\|_\infty > \frac{C^*}{2}\sqrt{\frac{2^{ld}ld}{n}} \right) \end{array} \right\}$$

$$\leq \sum_{l>l^*} 2\exp\left(-Cld\right),$$

$$(6.9)$$

In the fourth and fifth of the above series of inequalities, we have used (6.3) and the definition of $l^*$ respectively. The last inequality in the above display holds for a $C > 0$ depending on $B_U, \psi_{0,0}^0, \psi_{0,0}^1$ and the inequality follows from Lemma B.6 provided we choose $C^*$ large enough depending on $M, B_U, \psi_{0,0}^0, \psi_{0,0}^1, \gamma_{\max}$. In particular, this implies that, choosing $C^*$ large enough will guarantee that there exists a $\eta > 3$ such that for large enough $n$, one has for any $l > l^*$

$$\mathbb{P}(\hat{l} = l) \leq n^{-\eta}. \qquad (6.10)$$

This along with 6.8 and choice of $l_{\max}$ implies that

$$II \leq C_3 \sum_{l>l^*} 2^{ld} n^{-\eta} = C_3 \sum_{l>l^*} \frac{2^{ld}}{n} n^{-\eta+1} \leq \frac{l_{\max}}{n^{\eta-1}} \leq \frac{\log n}{n}. \qquad (6.11)$$

Finally combining equations (6.7) and (6.11), we have the existence of an estimator $\widetilde{g}$ depending on $M, B_U$, and $\gamma_{\max}$ (once we have fixed our choice of father and mother wavelets), such that for every $(\beta, \gamma) \in [\beta_{\min}, \beta_{\max}] \times [\gamma_{\min}, \gamma_{\max}]$,

$$\sup_{P \in \mathcal{P}(\beta,\gamma)} \mathrm{E}_P \|\widetilde{g} - g\|_\infty \leq (C)^{\frac{d}{2\gamma+d}} \left(\frac{n}{\log n}\right)^{-\frac{\gamma}{2\gamma+d}},$$

with a large enough positive $C$ depending on $M, B_U$, and $\gamma_{\max}$.

We next show that uniformly over $P \in \mathcal{P}(\beta, \gamma)$, $\widetilde{g}$ belongs to $H(\gamma, C)$ with probability at least $1 - 1/n^2$, for a large enough constant $C$ depending on $M, B_U$, and $\gamma_{\max}$. Towards this end, note that, for any $C > 0$ and $l' > 0$, (letting for any $h \in L_2[0,1]^d$, $\|\langle h, \psi_{l',.}\rangle\|_2$ be the vector $L_2$ norm of the vector $\left(\langle h, \psi_{l',k'}^v\rangle : k' \in \mathcal{Z}_{l'}, v \in \{0,1\}^d - \{0\}^d\right)$). We have,

$$\mathbb{P}_{P,2}\left( 2^{l'(\gamma+\frac{d}{2})} \|\langle \widetilde{g}, \psi_{l',.}\rangle\|_\infty > C \right)$$

$$= \sum_{l=l_{\min}}^{l_{\max}} \mathbb{P}_{P,2}\left( 2^{l'(\gamma+\frac{d}{2})} \|\langle \hat{g}_l, \psi_{l',.}\rangle\|_\infty > C, \hat{l} = l \right) \mathcal{I}\left(l' \leq l\right)$$

$$= \sum_{l=l_{\min}}^{l^*} \mathbb{P}_{P,2}\left( 2^{l'(\gamma+\frac{d}{2})} \|\langle \hat{g}_l, \psi_{l',.}\rangle\|_\infty > C, \hat{l} = l \right) \mathcal{I}\left(l' \leq l\right)$$

$$+ \sum_{l=l^*+1}^{l_{\max}} \mathbb{P}_{P,2}\left(2^{l'(\gamma+\frac{d}{2})}\|\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle\|_\infty > C, \hat{l} = l\right)\mathcal{I}\left(l' \leq l\right)$$

$$\leq \sum_{l=l_{\min}}^{l^*} \mathbb{P}_{P,2}\left(2^{l'(\gamma+\frac{d}{2})}\|\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle\|_\infty > C\right)\mathcal{I}\left(l' \leq l\right)$$

$$+ \sum_{l=l^*+1}^{l_{\max}} \mathbb{P}_{P,2}\left(\hat{l} = l\right)\mathcal{I}\left(l' \leq l\right)$$

$$\leq \sum_{l=l_{\min}}^{l^*} \mathbb{P}_{P,2}\left(2^{l'(\gamma+\frac{d}{2})}\|\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle\|_\infty > C\right)\mathcal{I}\left(l' \leq l\right) + \sum_{l>l^*} n^{-\eta}, \tag{6.12}$$

where the last inequality follows from (6.10) for some $\eta > 3$ provided $C^*$ is chosen large enough as before. Now,

$$\mathbb{P}_{P,2}\left(2^{l'(\gamma+\frac{d}{2})}\|\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle\|_\infty > C\right)$$

$$\leq \mathbb{P}_{P,2}\left(2^{l'(\gamma+\frac{d}{2})}\|\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle - \mathrm{E}_{P,2}\left(\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle\right)\|_\infty > C/2\right)$$

$$+ \mathbb{P}_{P,2}\left(2^{l'(\gamma+\frac{d}{2})}\|\mathrm{E}_{P,2}\left(\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle\right)\|_\infty > C/2\right)$$

$$= \mathbb{P}_{P,2}\left(2^{l'(\gamma+\frac{d}{2})}\|\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle - \mathrm{E}_{P,2}\left(\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle\right)\|_\infty > C/2\right)$$

if $C > 2M$ (by definition 5.5). Therefore, from (6.12), one has for any $C > 2M$,

$$\mathbb{P}_{P,2}\left(2^{l'(\gamma+\frac{d}{2})}\|\langle \hat{g}, \boldsymbol{\psi}_{l',.}\rangle\|_\infty > C\right)$$

$$\leq \sum_{l=l_{\min}}^{l^*} \mathbb{P}_{P,2}\left(2^{l'(\gamma+\frac{d}{2})}\|\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle - \mathrm{E}_{P,2}\left(\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle\right)\|_\infty > C/2\right)\mathcal{I}\left(l' \leq l\right)$$

$$+ \sum_{l=l^*}^{l_{\max}} n^{-3}\mathcal{I}\left(l' \leq l\right). \tag{6.13}$$

Considering the first term of the last summand of the above display, we have

$$\sum_{l=l_{\min}}^{l^*} \mathbb{P}_{P,2}\left(2^{l'(\gamma+\frac{d}{2})}\|\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle - \mathrm{E}_{P,2}\left(\langle \hat{g}_l, \boldsymbol{\psi}_{l',.}\rangle\right)\|_\infty > C/2\right)\mathcal{I}\left(l' \leq l\right)$$

$$= \sum_{l=l_{\min}}^{l^*} \sum_{k\in\mathcal{Z}_{l'}} \sum_{v\in\{0,1\}^d} \mathbb{P}_{P,2}\left(\left|\frac{1}{n}\sum_{i=n+1}^{2n}\left(\psi_{l',k}^v(\mathbf{X}_i) - \mathrm{E}_{P,2}\left(\psi_{l',k}^v(\mathbf{X}_i)\right)\right)\right| > \frac{C/2}{2^{l'(\gamma+\frac{d}{2})}}\right)\mathcal{I}\left(l' \leq l\right)$$

By Bersntein's Inequality, for any $\lambda > 0$,

$$\mathbb{P}_{P,2}\left(\left|\frac{1}{n}\sum_{i=n+1}^{2n}\left(\psi_{l',k}^v(\mathbf{X}_i) - \mathrm{E}_{P,2}\left(\psi_{l',k}^v(\mathbf{X}_i)\right)\right)\right| > \lambda\right)$$

$$\leq 2\exp\left(-\frac{n\lambda^2}{2\left(\sigma^2 + \|\psi_{l',k}^v\|_\infty\lambda/3\right)}\right),$$

where $\sigma^2 = \mathrm{E}_{P,2}\left(\psi^v_{l',k}(\mathbf{X}_i) - \mathrm{E}_{P,2}\left(\psi^v_{l',k}(\mathbf{X}_i)\right)\right)^2$. Indeed, there exists constant $C_4$ depending on the $\psi^0_{0,0}, \psi^1_{0,0}, \gamma_{\max}$ such that $\sigma^2 \leq C_4$ and $\|\psi^v_{l',k}\|_\infty \leq C_4 2^{\frac{l'd}{2}}$. Therefore,

$$\sum_{l=l_{\min}}^{l^*} \sum_{k \in \mathcal{Z}_{l'}} \sum_{v \in \{0,1\}^d} \mathbb{P}_{P,2}\left(\left|\frac{1}{n}\sum_{i=n+1}^{2n}\left(\psi^v_{l',k}(\mathbf{X}_i) - \mathrm{E}_{P,2}\left(\psi^v_{l',k}(\mathbf{X}_i)\right)\right)\right| > \frac{C/2}{2^{l'(\gamma+\frac{d}{2})}}\right)\mathcal{I}\left(l' \leq l\right)$$

$$\leq 2\sum_{l=l_{\min}}^{l^*} \sum_{k \in \mathcal{Z}_{l'}} \sum_{v \in \{0,1\}^d} \exp\left(-\frac{C^2}{8C_4}\frac{n2^{-2l'(\gamma+\frac{d}{2})}}{1+\frac{C}{2}2^{\frac{l'd}{2}}2^{-l'(\gamma+\frac{d}{2})}}\right)\mathcal{I}\left(l' \leq l\right)$$

$$= 2\sum_{l=l_{\min}}^{l^*} \sum_{k \in \mathcal{Z}_{l'}} \sum_{v \in \{0,1\}^d} \exp\left(-\frac{C^2}{8C_4}\frac{n2^{-2l'\gamma}}{2^{l'd}+\frac{C}{2}2^{l'(d-\gamma)}}\right)\mathcal{I}\left(l' \leq l\right)$$

$$\leq 2\sum_{l=l_{\min}}^{l^*} \sum_{k \in \mathcal{Z}_{l'}} \sum_{v \in \{0,1\}^d} \exp\left(-\frac{C^2}{8(1+\frac{C}{2})C_4}\frac{n2^{-2l'\gamma}}{2^{l'd}}\right)\mathcal{I}\left(l' \leq l\right)$$

$$= 2\sum_{l=l_{\min}}^{l^*} \sum_{k \in \mathcal{Z}_{l'}} \sum_{v \in \{0,1\}^d} \exp\left(-\frac{C^2}{8(1+\frac{C}{2})C_4}\frac{n2^{-2l^*\gamma}}{2^{l^*d}l^*d}2^{(l^*-l)(d+2\gamma)}l^*d\right)\mathcal{I}\left(l' \leq l\right)$$

$$\leq 2\sum_{l=l_{\min}}^{l^*} \sum_{k \in \mathcal{Z}_{l'}} \sum_{v \in \{0,1\}^d} \exp\left(-\frac{C^2C_2}{2^{d+3}C_1(1+\frac{C}{2})C_4}2^{(l^*-l)(d+2\gamma)}l^*d\right)\mathcal{I}\left(l' \leq l\right) \quad (6.14)$$

$$\leq 2\sum_{l=l_{\min}}^{l^*} \sum_{k \in \mathcal{Z}_{l'}} \sum_{v \in \{0,1\}^d} \exp\left(-\frac{C^2C_2}{2^{d+3}C_1(1+\frac{C}{2})C_4}ld\right)\mathcal{I}\left(l' \leq l\right)$$

$$\leq 2\sum_{l=l_{\min}}^{l^*} C(\psi^0_{0,0},\psi^1_{0,0})2^{l'd}\exp\left(-\frac{C^2C_2}{2^{d+3}C_1(1+\frac{C}{2})C_4}ld\right)\mathcal{I}\left(l' \leq l\right)$$

$$\leq 2\sum_{l=l_{\min}}^{l^*} C(\psi^0_{0,0},\psi^1_{0,0})\exp\left(-\left(\frac{C^2C_2}{2^{d+3}C_1(1+\frac{C}{2})C_4}-1\right)ld\right)\mathcal{I}\left(l' \leq l\right)$$

$$\leq 2l_{\max}C(\psi^0_{0,0},\psi^1_{0,0})\exp\left(-\left(\frac{C^2C_2}{2^{d+3}C_1(1+\frac{C}{2})C_4}-1\right)l_{\min}d\right)\mathcal{I}\left(l' \leq l_{\max}\right),$$

if $\frac{C^2C_2}{2^{d+3}C_1(1+\frac{C}{2})C_4} \geq 1$. Above inequality 6.14 uses the definition of $l^*$. Indeed choosing $C$ large enough, one can guarantee, $\left(\frac{C^2C_2}{2^{d+3}C_1(1+\frac{C}{2})C_4}-1\right)l_{\min}d \geq 4\log n$. Such a choice of $C$ implies that,

$$\sum_{l=l_{\min}}^{l^*} \sum_{k \in \mathcal{Z}_{l'}} \sum_{v \in \{0,1\}^d} \mathbb{P}_{P,2}\left(\left|\frac{1}{n}\sum_{i=n+1}^{2n}\left(\psi^v_{l',k}(\mathbf{X}_i) - \mathrm{E}_{P,2}\left(\psi^v_{l',k}(\mathbf{X}_i)\right)\right)\right| > \frac{C/2}{2^{l'(\gamma+\frac{d}{2})}}\right)\mathcal{I}\left(l' \leq l\right)$$

$$\leq \frac{C(\psi^0_{0,0},\psi^1_{0,0})}{n^3}\mathcal{I}(l' \leq l_{\max}),$$

which in turn implies that, for $C$ sufficiently large (depending on $M, \psi^0_{0,0}, \psi^1_{0,0}$) one has

$$\mathbb{P}_{P,2}\left(2^{l'(\gamma+\frac{d}{2})}\|\langle \hat{g}, \psi_{l',.}\rangle\|_2 > C\right) \leq \frac{C(\psi^0_{0,0},\psi^1_{0,0})+1}{n^3}\mathcal{I}(l' \leq l_{\max}).$$

This along with the logarithmic in $n$ size of $l_{\max}$ implies that for sufficiently large $n$, uniformly over $P \in \mathcal{P}(\beta, \gamma)$, $\widetilde{g}$ belongs to $H(\gamma, C)$ with probability at least $1 - 1/n^2$, for a large enough constant $C$ depending on $M, B_U$, and $\gamma_{\max}$ (the choice of $\psi_{0,0}^0, \psi_{0,0}^1$ being fixed by specifying a regularity $S > \gamma_{\max}$).

However this $\widetilde{g}$ does not satisfy the desired point-wise bounds. To achieve this let $\phi$ be a $C^\infty$ function such that $\psi(x)|_{[B_L, B_U]} \equiv x$ while $\frac{B_L}{2} \leq \psi(\mathbf{x}) \leq 2B_U$ for all $x$. Finally, consider the estimator $\hat{g}(\mathbf{x}) = \psi(\widetilde{g}(\mathbf{x}))$. We note that $|g(\mathbf{x}) - \hat{g}(\mathbf{x})| \leq |g(\mathbf{x}) - \widetilde{g}(\mathbf{x})|$— thus $\hat{g}$ is adaptive to the smoothness of the design density. The boundedness of the constructed estimator follows from the construction. Finally, we wish to show that almost surely, the constructed estimator belongs to the Hölder space with the same smoothness, possibly of a different radius. This is captured by the next lemma, proof of which can be completed by following arguments similar to proof of Lemma 3.1 in Mukherjee and Sen (2016). In particular,

LEMMA 6.2. *For all $h \in H(\beta, M)$, $\psi(h) \in H(\beta, C(M, \beta))$, where $C(M, \beta)$ is a universal constant dependent only on $M, \beta$ and independent of $h \in H(\beta, M)$.*

Now, the construction of $\hat{f}$ satisfying the desired properties of Theorem 2.3 can be done following ideas from proof of Theorem 1.1 of Mukherjee and Sen (2016). In particular, construct the estimator $\hat{g}$ of the design density $g$ as above from second part of the sample and let for $j \in \mathcal{T}_1$, $\hat{f}_j(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \frac{W_i}{\hat{g}(\mathbf{X}_i)} K_{V_j}(\mathbf{X}_i, \mathbf{x})$. Now, let

$$\hat{j} = \min \left\{ j \in \mathcal{T}_1 : \|\hat{f}_j - \hat{f}_{j'}\|_\infty \leq C^{**} \sqrt{\frac{2^{j'd} j'd}{n}}, \ \forall j' \in \mathcal{T}_1 \text{ s.t. } j' \geq j \right\}.$$

where $C^{**}$ depends only on the known parameters of the problem and can be determined from the proof hereafter. Thereafter, consider the estimator $\widetilde{f} := \hat{f}_{\hat{j}}$.

Now define

$$j^* := \min \left\{ j \in \mathcal{T}_1 : 2^{-jd\frac{\beta}{d}} \leq \sqrt{\frac{2^{jd} jd}{n}} \right\},$$

Therefore

$$\mathrm{E}_P \|\widetilde{f} - f\|_\infty \leq \mathrm{E}_P \|\hat{f}_{\hat{j}} - f\|_\infty \mathcal{I}(\hat{j} \leq j^*) + \mathrm{E}_P \|\hat{f}_{\hat{j}} - f\|_\infty \mathcal{I}(\hat{j} > j^*). \tag{6.15}$$

Thereafter using Lemma B.6 and 6.3 we have

$$\begin{aligned}
&\mathrm{E}_P \|\hat{f}_{\hat{j}} - f\|_\infty \mathcal{I}(\hat{j} \leq j^*) \\
&\leq \mathrm{E}_P \|\hat{f}_{\hat{j}} - \hat{f}_{j^*}\|_\infty \mathcal{I}(\hat{j} \leq j^*) + \mathrm{E}_P \|\hat{f}_{j^*} - f\|_\infty \\
&\leq C^{**} \sqrt{\frac{2^{j^*d} j^*d}{n}} + \mathrm{E}_{P,2} \|\hat{f}_{j^*} - \mathrm{E}_{P,1}(\hat{f}_{j^*})\|_\infty + \mathrm{E}_{P,2} \|\mathrm{E}_{P,1}(\hat{f}_{j^*}) - f\|_\infty \\
&\leq (C^{**} + C(B_U, B_L, \psi_{0,0}^0, \psi_{0,0}^1)) \sqrt{\frac{2^{j^*d} j^*d}{n}} \\
&\quad + \mathrm{E}_{P,2} \|\Pi(f(\frac{g}{\hat{g}} - 1)|V_{j^*})\|_\infty + \|f - \Pi(f|V_{j^*})\|_\infty \\
&\leq (C^{**} + C(B_U, B_L, \psi_{0,0}^0, \psi_{0,0}^1)) \sqrt{\frac{2^{j^*d} j^*d}{n}}
\end{aligned}$$

$$+ C(M, \psi_{0,0}^0, \psi_{0,0}^1) 2^{-j^*\beta} + \mathrm{E}_{P,2} \|\Pi(f(\tfrac{g}{\hat{g}} - 1)|V_{j^*})\|_\infty.$$

$$(6.16)$$

Now, by standard computations involving compactly wavelet bases and property of $\hat{g}$

$$
\begin{aligned}
\mathrm{E}_{P,2}\|\Pi(f(\tfrac{g}{\hat{g}} - 1)|V_{j^*})\|_\infty &\le C(\psi_{0,0}^0, \psi_{0,0}^1)\mathrm{E}_{P,2}\|f(\tfrac{g}{\hat{g}} - 1)\|_\infty \\
&\le C(B_U, B_L, \psi_{0,0}^0, \psi_{0,0}^1)\mathrm{E}_{P,2}\|\hat{g} - g\|_\infty \\
&\le C(B_U, B_L, M, \gamma_{\max}, \psi_{0,0}^0, \psi_{0,0}^1)\left(\frac{n}{\log n}\right)^{-\frac{\gamma}{2\gamma+d}}.
\end{aligned}
$$

$$(6.17)$$

Combining (6.16), (6.17), definition of $j^*$, and the fact that $\gamma > \beta$, we have

$$\mathrm{E}_P\|\hat{f}_{\hat{j}} - f\|_\infty \mathcal{I}(\hat{j} \le j^*) \le C(B_U, B_L, M, \gamma_{\max}, \psi_{0,0}^0, \psi_{0,0}^1)\left(\frac{n}{\log n}\right)^{-\frac{\beta}{2\beta+d}}. \qquad (6.18)$$

provided $C^* *$ is chosen depending only the known parameters of the problem. Now using arguments similar to those leading to (6.8) we have

$$\mathrm{E}_P\|\hat{f}_{\hat{j}} - f\|_\infty \mathcal{I}(\hat{j} > j^*) \le C(B_U, B_L, \psi_{0,0}^0, \psi_{0,0}^1)\sum_{j > j^*} 2^{jd}\mathbb{P}_P(\hat{j} = j). \qquad (6.19)$$

We now complete the control over II by suitably bounding $\mathbb{P}_P(\hat{j} = j)$. To this end, note that for any $j > j^*$,

$$
\begin{aligned}
&\mathbb{P}_P(\hat{j} = j) \\
&\le \sum_{j > j^*} \mathbb{P}_P\left(\|\hat{f}_j - \hat{f}_{j^*}\|_\infty > C^{**}\sqrt{\frac{2^{jd}jd}{n}}\right) \\
&\le \sum_{j > j^*} \mathrm{E}_{P,2}\left\{
\begin{array}{c}
\mathbb{P}_{P,1}\left(\|\hat{f}_{j^*} - \mathrm{E}_{P,1}\left(\hat{f}_{j^*}\right)\|_\infty > \frac{C^{**}}{2}\sqrt{\frac{2^{jd}jd}{n}} - \|\mathrm{E}_{P,1}\left(\hat{f}_{j^*}\right) - \mathrm{E}_{P,1}\left(\hat{f}_j\right)\|_\infty\right) \\
+\mathbb{P}_{P,1}\left(\|\hat{f}_j - \mathrm{E}_{P,1}\left(\hat{f}_j\right)\|_\infty > \frac{C^{**}}{2}\sqrt{\frac{2^{jd}jd}{n}}\right)
\end{array}
\right\} \\
&\le \sum_{j > j^*} \mathrm{E}_{P,2}\left\{
\begin{array}{c}
\mathbb{P}_{P,1}\left(\|\hat{f}_{j^*} - \mathrm{E}_{P,1}\left(\hat{f}_{j^*}\right)\|_\infty > \frac{C^{**}}{2}\sqrt{\frac{2^{jd}jd}{n}} - \|\Pi\left(f\tfrac{g}{\hat{g}}|V_{j^*}\right) - \Pi\left(f\tfrac{g}{\hat{g}}|V_j\right)\|_\infty\right) \\
+\mathbb{P}_{P,1}\left(\|\hat{f}_j - \mathrm{E}_{P,1}\left(\hat{f}_j\right)\|_\infty > \frac{C^{**}}{2}\sqrt{\frac{2^{jd}jd}{n}}\right)
\end{array}
\right\}.
\end{aligned}
$$

Now,

$$
\begin{aligned}
&\left\|\Pi\left(f\tfrac{g}{\hat{g}}|V_{j^*}\right) - \Pi\left(f\tfrac{g}{\hat{g}}|V_j\right)\right\|_\infty \\
&\le C(M, \psi_{0,0}^0, \psi_{0,0}^1)2^{-j^*\beta} + C(B_U, B_L, , \psi_{0,0}^0, \psi_{0,0}^1)\|\hat{g} - g\|_\infty.
\end{aligned}
$$

Using the fact that $\sqrt{\frac{2^{jd}jd}{n}} > \sqrt{\frac{2^{j^*d}j^*d}{n}}$ for $j > j^*$, we have using the definition of $j^*$ that there exists $C, C' > 0$ depending on $M, B_U, B_L, \psi_{0,0}^0, \psi_{0,0}^1$ such that

$$\mathbb{P}_P(\hat{j} = j)$$

$$\leq \sum_{j>j^*} E_{P,2} \left\{ \begin{array}{c} \mathbb{P}_{P,1}\left( \|\hat{f}_{j^*} - \mathrm{E}_{P,1}\left(\hat{f}_{j^*}\right)\|_\infty > (\frac{C^{**}}{2} - C)\sqrt{\frac{2^{jd}jd}{n}} \right) \\ +\mathbb{P}_{P,1}\left( \|\hat{f}_j - \mathrm{E}_{P,1}\left(\hat{f}_j\right)\|_\infty > \frac{C^{**}}{2}\sqrt{\frac{2^{jd}jd}{n}} \right) + \mathbb{P}_{P,2}\left( \|\hat{g} - g\|_\infty > C'\sqrt{\frac{2^{j^*d}j^*d}{n}} \right) \end{array} \right\}.$$

(6.20)

Now, provided $C^{**} > 2C$ is chosen large enough (depending on $B_U, B_L, \psi_{0,0}^0, \psi_{0,0}^1$) we have there exists large enough $C''$ (depending on $B_U, B_L, \psi_{0,0}^0, \psi_{0,0}^1$) such that

$$\mathbb{P}_{P,1}\left( \|\hat{f}_{j^*} - \mathrm{E}_{P,1}\left(\hat{f}_{j^*}\right)\|_\infty > (\frac{C^{**}}{2} - C)\sqrt{\frac{2^{jd}jd}{n}} \right)$$
$$+ \mathbb{P}_{P,1}\left( \|\hat{f}_j - \mathrm{E}_{P,1}\left(\hat{f}_j\right)\|_\infty > \frac{C^{**}}{2}\sqrt{\frac{2^{jd}jd}{n}} \right) \leq 2e^{-C''jd}.$$

(6.21)

Henceforth, whenever required, $C, C', C''$ will be chosen to be large enough depending on the known parameters of the problem, which in turn will imply that $C^{**}$ can be chosen large enough depending on the known parameters of the problem as well. First note that, the last term in the above display can be bounded rather crudely using the following lemma.

LEMMA 6.3. *Assume $\gamma_{\min} > \beta_{\max}$. Then for $C', C_1, C_2 > 0$ (chosen large enough depending on $B_U, \psi_{0,0}^0, \psi_{0,0}^1$) one has*

$$\sup_{P\in\mathcal{P}(\beta,\gamma)} \mathbb{P}_{P,2}\left( \|\hat{g} - g\|_\infty > C'\sqrt{\frac{2^{j^*d}j^*d}{n}} \right) \leq C_1(l_{\max} - l_{\min})e^{-C_2 l_{\min}d}.$$

The proof of Lemma 6.3 can be argued as follows. Indeed, $\hat{g} = \psi(\widetilde{g})$, where $\psi(x)$ is $C^\infty$ function which is identically equal to $x$ on $[B_L, B_U]$ and has universally bounded first derivative. Therefore, it is enough to prove Lemma 6.3 for $\widetilde{g}$ instead of $\hat{g}$ and thereby invoking a simple first order Taylor series argument along with the fact that $\psi(g) \equiv g$ owing to the bounds on $g$. The crux of the argument for proving Lemma 6.3 is that by Lemma B.6, any $\hat{g}_l$ for $l \in \mathcal{T}_2$ suitably concentrates around $g$ in a radius of the order of $\sqrt{\frac{2^{ld}ld}{n}}$. The proof of the lemma is therefore very similar to the proof of adaptivity of $\hat{g}$ (by dividing into cases where the chosen $\hat{l}$ is larger and smaller than $l^*$ respectively and thereafter invoking Lemma B.6) and therefore we omit the details.

Plugging in the result of Lemma 6.3 into (6.20), and thereafter using the facts that $\gamma_{\min} > \beta_{\max}$, $l_{\max}, j_{\max}$ are both poly logarithmic in nature, along with equations (6.15), and (6.18), (6.19), (6.21) we have the existence of an estimator $\widetilde{f}$ depending on $M, B_U, B_L, \beta_{\min}, \beta_{\max}, \gamma_{\max}$, such that for every $(\beta, \gamma) \in [\beta_{\min}, \beta_{\max}] \times [\gamma_{\min}, \gamma_{\max}]$,

$$\sup_{P\in\mathcal{P}(\beta,\gamma)} \mathrm{E}_P\|\widetilde{f} - f\|_\infty \leq C\left(\frac{n}{\log n}\right)^{-\frac{\beta}{2\beta+d}},$$

with a large enough positive constant $C$ depending on $M, B_U, B_L, \beta_{\min}, \gamma_{\max}, \psi_{0,0}^0, \psi_{0,0}^1$.

However this $\widetilde{f}$ does not satisfy the desired point-wise bounds. To achieve this, as before, let $\phi$ be a $C^\infty$ function such that $\psi(x)|_{[B_L, B_U]} \equiv x$ while $\frac{B_L}{2} \leq \psi(\mathbf{x}) \leq 2B_U$ for all $x$. Finally, consider the estimator $\hat{f}(\mathbf{x}) = \psi(\widetilde{g}(\mathbf{x}))$. We note that $|f(\mathbf{x}) - \hat{f}(\mathbf{x})| \leq |f(\mathbf{x}) - \widetilde{f}(\mathbf{x})|$— thus $\hat{f}$ is adaptive to the smoothness of the design density. The boundedness of the constructed estimator follows from

the construction. Finally,the proof of the fact that the constructed estimator belongs to the Hölder space with the same smoothness, possibly of a different radius follows once again from of Lemma 6.2.

■

## References.

Baraud, Y. (2002). Model selection for regression on a random design. *ESAIM: Probability and Statistics* **6** 127–146.

Bickel, P. J. and Ritov, Y. (1988). Estimating integrated squared density derivatives: sharp best order of convergence estimates. *Sankhyā: The Indian Journal of Statistics, Series A* 381–393.

Bickel, P. J., Klaassen, C. A., Bickel, P. J., Ritov, Y., Klaassen, J., Wellner, J. A. and Ritov, Y. (1993). *Efficient and adaptive estimation for semiparametric models*. Johns Hopkins University Press Baltimore.

Birgé, L. and Massart, P. (1995). Estimation of integral functionals of a density. *The Annals of Statistics* 11–29.

Brown, L. D. and Levine, M. (2007). Variance estimation in nonparametric regression via the difference sequence method. *The Annals of Statistics* **35** 2219–2232.

Brown, L. D. and Low, M. G. (1996). A constrained risk inequality with applications to nonparametric functional estimation. *The annals of Statistics* **24** 2524–2535.

Bull, A. D. and Nickl, R. (2013). Adaptive confidence sets in Lˆ 2. *Probability Theory and Related Fields* **156** 889–919.

Cai, T. T. and Low, M. G. (2003). A note on nonparametric estimation of linear functionals. *Annals of statistics* 1140–1153.

Cai, T. T. and Low, M. G. (2004). Minimax estimation of linear functionals over nonconvex parameter spaces. *Annals of statistics* 552–576.

Cai, T. T. and Low, M. G. (2005a). Nonquadratic estimators of a quadratic functional. *The Annals of Statistics* 2930–2956.

Cai, T. T. and Low, M. G. (2005b). On adaptive estimation of linear functionals. *The Annals of Statistics* **33** 2311–2343.

Cai, T. T. and Low, M. G. (2006). Optimal adaptive estimation of a quadratic functional. *The Annals of Statistics* **34** 2298–2325.

Cai, T. T. and Low, M. G. (2011). Testing composite hypotheses, Hermite polynomials and optimal estimation of a nonsmooth functional. *The Annals of Statistics* **39** 1012–1041.

Cai, T. T. and Wang, L. (2008). Adaptive variance function estimation in heteroscedastic nonparametric regression. *The Annals of Statistics* **36** 2025–2054.

Cohen, A., Daubechies, I. and Vial, P. (1993). Wavelets on the interval and fast wavelet transforms. *Applied and computational harmonic analysis* **1** 54–81.

Crump, R. K., Hotz, V. J., Imbens, G. W. and Mitnik, O. A. (2009). Dealing with limited overlap in estimation of average treatment effects. *Biometrika* asn055.

Donoho, D. L., Liu, R. C. and MacGibbon, B. (1990). Minimax risk over hyperrectangles, and implications. *The Annals of Statistics* 1416–1437.

Donoho, D. L. and Nussbaum, M. (1990). Minimax quadratic estimation of a quadratic functional. *Journal of Complexity* **6** 290–323.

Efromovich, S. and Low, M. G. (1994). Adaptive estimates of linear functionals. *Probability theory and related fields* **98** 261–275.

Efromovich, S. and Low, M. (1996). On optimal adaptive estimation of a quadratic functional. *The Annals of Statistics* **24** 1106–1125.

Efromovich, S. and Samarov, A. (2000). Adaptive estimation of the integral of squared regression derivatives. *Scandinavian journal of statistics* **27** 335–351.

Fan, J. (1991). On the estimation of quadratic functionals. *The Annals of Statistics* 1273–1294.

Fan, J. and Yao, Q. (1998). Efficient estimation of conditional variance functions in stochastic regression. *Biometrika* **85** 645–660.

Giné, E., Latala, R. and Zinn, J. (2000). Exponential and moment inequalities for U-statistics. In *High Dimensional Probability II* 13–38. Springer.

Giné, E. and Nickl, R. (2008). A simple adaptive estimator of the integrated square of a density. *Bernoulli* 47–61.

Giné, E. and Nickl, R. (2015). Mathematical foundations of infinite-dimensional statistical models. *Cambridge Series in Statistical and Probabilistic Mathematics*.

Hall, P. and Carroll, R. (1989). Variance function estimation in regression: the effect of estimating the mean. *Journal of the Royal Statistical Society. Series B (Methodological)* 3–14.

Hall, P. and Marron, J. S. (1987). Estimation of integrated squared density derivatives. *Statistics & Probability Letters* **6** 109–115.

HÄRDLE, W., KERKYACHARIAN, G., TSYBAKOV, A. and PICARD, D. (1998). *Wavelets, approximation, and statistical applications.* Springer.

HOUDRÉ, C. and REYNAUD-BOURET, P. (2003). Exponential inequalities, with constants, for U-statistics of order two. 55–69.

IBRAGIMOV, I. A. and HAS' MINSKII, R. Z. (2013). *Statistical estimation: asymptotic theory* **16**. Springer Science & Business Media.

KERKYACHARIAN, G. and PICARD, D. (1996). Estimating nonquadratic functionals of a density using Haar wavelets. *The Annals of Statistics* **24** 485–507.

KLEMELA, J. and TSYBAKOV, A. B. (2001). Sharp adaptive estimation of linear functionals. *The Annals of statistics* 1567–1600.

LAURENT, B. (1996). Efficient estimation of integral functionals of a density. *The Annals of Statistics* **24** 659–681.

LAURENT, B. and MASSART, P. (2000). Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics* 1302–1338.

LEPSKI, O. (1991). On a problem of adaptive estimation in Gaussian white noise. *Theory of Probability & Its Applications* **35** 454–466.

LEPSKI, O. V. (1992). On problems of adaptive estimation in white Gaussian noise. *Topics in nonparametric estimation* **12** 87–106.

LOW, M. G. (1992). Renormalization and white noise approximation for nonparametric functional estimation problems. *The Annals of Statistics* 545–554.

MUKHERJEE, R., NEWEY, W. K. and ROBINS, J. M. (2017). Semiparametric Efficient Empirical Higher Order Influence Function Estimators. *arXiv preprint arXiv:1705.07577.*

MUKHERJEE, R. and SEN, S. (2016). Optimal Adaptive Inference in Random Design Binary Regression. *Bernoulli (To Appear).*

MUKHERJEE, R., TCHETGEN TCHETGEN, E. and ROBINS, J. (2017). Supplement to "Adpative Estimation of Nonparametric Functionals".

NEMIROVSKI, A. (2000). Topics in non-parametric. *Ecole dEté de Probabilités de Saint-Flour* **28** 85.

PETROV, V. V. (1995). Limit Theorems of Probability Theory. Sequences of Independent Random Variables, vol. 4 of. *Oxford Studies in Probability.*

ROBINS, J. M., MARK, S. D. and NEWEY, W. K. (1992). Estimating exposure effects by modelling the expectation of exposure conditional on confounders. *Biometrics* 479–495.

ROBINS, J., LI, L., TCHETGEN, E. and VAN DER VAART, A. (2008). Higher order influence functions and minimax estimation of nonlinear functionals. In *Probability and Statistics: Essays in Honor of David A. Freedman* 335–421. Institute of Mathematical Statistics.

ROBINS, J., TCHETGEN, E. T., LI, L. and VAN DER VAART, A. (2009). Semiparametric minimax rates. *Electronic Journal of Statistics* **3** 1305–1321.

ROBINS, J., LI, L., MUKHERJEE, R., TCHETGEN, E. T. and VAN DER VAART, A. (2016). Higher Order Estimating Equations for High-dimensional Models. *The Annals of Statistics (To Appear).*

RUPPERT, D., WAND, M. P., HOLST, U. and HÖSJER, O. (1997). Local polynomial variance-function estimation. *Technometrics* **39** 262–273.

SHAO, Q.-M. (2000). A comparison theorem on moment inequalities between negatively associated and independent random variables. *Journal of Theoretical Probability* **13** 343–356.

TRIBOULEY, K. (2000). Adaptive estimation of integrated functionals. *Mathematical Methods of Statistics* **9** 19–38.

TSIATIS, A. (2007). *Semiparametric theory and missing data.* Springer Science & Business Media.

VAN DER VAART, A. W. (2000). *Asymptotic statistics* **3**. Cambridge university press.

## APPENDIX A: PROOF OF REMAINING THEOREMS

**Proof of Theorem 3.1.**

PROOF.
*(i) Proof of Upper Bound*

The general scheme of proof involves identifying a non-adaptive minimax estimator of $\phi(P)$ under the knowledge of $P \in \mathcal{P}_{(\alpha,\beta,\gamma)}$, demonstrating suitable bias and variance properties of this sequence of estimators, and thereafter invoking Theorem 2.1 to conclude. This routine can be carried out as follows. Without loss of generality assume that we have $3n$ samples $\{Y_i, A_i, \mathbf{X}_i\}_{i=1}^n$. Divide the samples in to 3 equal parts (with the $l^{\text{th}}$ part being indexed by $\{(l-1)n + 1, \ldots, ln\}$

for $l \in \{1,2,3\}$), estimate $g$ by $\hat{g}$ adaptively from the third part (as in Theorem 2.3), and estimate $a$ and $b$ by $\hat{a}$ and $\hat{b}$ respectively, adaptively from the second part (as in Theorem 2.3). Let $E_{P,S}$ denote the expectation while samples with indices in $S$ held fixed, for $S \subset \{1,2,3\}$. A first order influence function for $\phi(P)$ at $P$ is given by $(Y - b(\mathbf{X})(A - a(\mathbf{X}))) - \phi(P)$ and a resulting first order estimator for $\phi(P)$ is $\frac{1}{n}\sum_{i=1}^{n}(Y_i - \hat{b}(\mathbf{X}_i))(A_i - \hat{a}(\mathbf{X}_i))$. This estimator has a bias $E_{P,\{2,3\}}\int\left((b(\mathbf{x}) - \hat{b}(\mathbf{x}))(a(\mathbf{x}) - \hat{a}(\mathbf{x}))\right)g(\mathbf{x})d\mathbf{x}$. Indeed for $\frac{\alpha+\beta}{2} < \frac{d}{2}$, this bias turns out to be sub-optimal compared to the minimax rate of convergence of $n^{-\frac{4\alpha+4\beta}{2\alpha+2\beta+d}}$ in mean squared loss. The most intuitive way to proceed is to estimate and correct for the bias. If there exists a "dirac-kernel" $K(\mathbf{x}_1,\mathbf{x}_2) \in L_2\left([0,1]^d \times [0,1]^d\right)$ such that $\int h(\mathbf{x}_1)K(\mathbf{x}_1,\mathbf{x}_2)dx_2 = h(\mathbf{x}_1)$ almost surely $\mathbf{x}_1$ for all $h \in L_2[0,1]^d$, then one can estimate the bias term by $\frac{1}{n(n-1)}\sum_{1 \le i_1 \ne i_2 \le n}\frac{(Y_{i_1}-\hat{b}(\mathbf{X}_{i_1})))}{\sqrt{g(\mathbf{X}_{i_1})}}K(\mathbf{X}_{i_1},\mathbf{X}_{i_2})\frac{(A_{i_2}-\hat{a}(\mathbf{X}_{i_2}))}{\sqrt{g(\mathbf{X}_{i_2})}}$, provided the marginal density $g$ was known. Indeed there are two concerns with the above suggestion. The first one being the knowledge of $g$. This can be relatively easy to deal with by plugging in an suitable estimate $\hat{g}$–although there are some subtleties involved (refer to to Section 4 for more on this). The primary concern though is the non-existence of a "dirac-kernel" of the above sort as an element of $L_2[0,1]^d \times L_2[0,1]^d$. This necessitates the following modification where one works with projection kernels on suitable finite dimensional linear subspace $L$ of $L_2[0,1]^d$ which guarantees existence of such kernels when the domain space is restricted to $L$. In particular, we work with the linear subspace $V_j$ (defined in 5) where the choice of $j$ is guided by the balance between the bias and variance properties of the resulting estimator. In particular, a choice of $2^j$ is guided by the knowledge of the parameter space $\mathcal{P}_{(\alpha,\beta,\gamma)}$. For any $j$ such that $n \le 2^{jd} \le n^2$, this implies that our bias corrected second order estimator of $\phi(P)$ is given by

$$\hat{\phi}_{n,j} = \frac{1}{n}\sum_{i=1}^{n}(Y_i - \hat{b}(\mathbf{X}_i))(A_i - \hat{a}(\mathbf{X}_i))$$
$$- \frac{1}{n(n-1)}\sum_{1 \le i_1 \ne i_2 \le n}S\left(\frac{(Y_{i_1} - \hat{b}(\mathbf{X}_{i_1})))}{\sqrt{\hat{g}(\mathbf{X}_{i_1})}}K_{V_j}(\mathbf{X}_{i_1},\mathbf{X}_{i_2})\frac{(A_{i_2} - \hat{a}(\mathbf{X}_{i_2}))}{\sqrt{\hat{g}(\mathbf{X}_{i_2})}}\right)$$

Note that division by $\hat{g}$ is permitted by the properties guaranteed by Theorem 2.3. Indeed this sequence of estimators is in the form of those considered by Theorem 2.1 with

$$L_1(O) = (Y - \hat{b}(\mathbf{X}))(A - \hat{a}(\mathbf{X})),$$

$$L_{2l}(O) = \frac{(Y - \hat{b}(\mathbf{X})))}{\sqrt{\hat{g}(\mathbf{X})}}, \quad L_{2r}(O) = \frac{(A - \hat{a}(\mathbf{X})))}{\sqrt{\hat{g}(\mathbf{X})}},$$

where by Theorem 2.3 $\max\{|L_1(O)|,|L_{2l}(O)|,|L_{2r}(O)|\} \le C(B_L,B_U)$. Therefore it remains to show that these sequence $\hat{\phi}_{n,j}$ satisfies the bias and variance property (A) and (B) necessary for application of Theorem 2.1.

We first verify the bias property. Utilizing the representation of the first order bias as stated above, we have

$$|E_P\left(\hat{\phi}_{n,j} - \phi(P)\right)|$$
$$= \left| \begin{array}{l} E_{P,\{2,3\}}\left[\int\left((b(\mathbf{x}) - \hat{b}(\mathbf{x}))(a(\mathbf{x}) - \hat{a}(\mathbf{x}))\right)g(\mathbf{x})d\mathbf{x}\right] \\ -E_P\left[S\left(\frac{(Y_1 - \hat{b}(\mathbf{X}_1)))}{\sqrt{\hat{g}(\mathbf{X}_1)}}K_{V_j}(\mathbf{X}_1,\mathbf{X}_2)\frac{(A_2 - \hat{a}(\mathbf{X}_2))}{\sqrt{\hat{g}(\mathbf{X}_2)}}\right)\right] \end{array} \right|$$

$$(A.1)$$

Now, using the notation $\delta b(\mathbf{x}) = b(\mathbf{x}) - \hat{b}(\mathbf{x})$ and $\delta a(\mathbf{x}) = a(\mathbf{x}) - \hat{a}(\mathbf{x})$, we have

$$
\mathrm{E}_P \left[ \frac{(Y_1 - \hat{b}(\mathbf{X}_1)))}{\sqrt{\hat{g}(\mathbf{X}_1)}} K_{V_j}(\mathbf{X}_1, \mathbf{X}_2) \frac{(A_2 - \hat{a}(\mathbf{X}_2))}{\sqrt{\hat{g}(\mathbf{X}_2)}} \right]
$$

$$
= E_{P,\{2,3\}} \int \int \left[ \frac{\delta b(\mathbf{x}_1) g(\mathbf{x}_1)}{\sqrt{\hat{g}(\mathbf{x}_1)}} K_{V_j}(\mathbf{x}_1, \mathbf{x}_2) \frac{\delta a(\mathbf{x}_2) g(\mathbf{x}_2)}{\sqrt{\hat{g}(\mathbf{x}_2)}} \right] d\mathbf{x}_1 d\mathbf{x}_2
$$

$$
= E_{P,\{2,3\}} \int \left[ \frac{\delta b(\mathbf{x}_1) g(\mathbf{x}_1)}{\sqrt{\hat{g}(\mathbf{x}_1)}} \Pi \left( \frac{\delta a g}{\sqrt{\hat{g}}} | V_j \right)(\mathbf{x}_1) \right] d\mathbf{x}_1
$$

$$
= E_{P,\{2,3\}} \int \left[ \frac{\delta b(\mathbf{x}_1) g(\mathbf{x}_1)}{\sqrt{\hat{g}(\mathbf{x}_1)}} \frac{\delta a(\mathbf{x}_1) g(\mathbf{x}_1)}{\sqrt{\hat{g}(\mathbf{x}_1)}} \right] d\mathbf{x}_1
$$

$$
- E_{P,\{2,3\}} \int \left[ \frac{\delta b(\mathbf{x}_1) g(\mathbf{x}_1)}{\sqrt{\hat{g}(\mathbf{x}_1)}} \Pi \left( \frac{\delta a g}{\sqrt{\hat{g}}} | V_j^{\perp} \right)(\mathbf{x}_1) \right] d\mathbf{x}_1
$$

$$
= \mathrm{E}_{P,\{2,3\}} \left[ \int \left( (b(\mathbf{x}) - \hat{b}(\mathbf{x}))(a(\mathbf{x}) - \hat{a}(\mathbf{x})) \right) g(\mathbf{x}) d\mathbf{x} \right]
$$

$$
+ E_{P,\{2,3\}} \int \left[ \delta a(\mathbf{x}_1) \delta b(\mathbf{x}_1) g^2(\mathbf{x}_1) \left( \frac{1}{\hat{g}(\mathbf{x}_1)} - \frac{1}{g(\mathbf{x}_1)} \right) \right] d\mathbf{x}_1
$$

$$
- E_{P,\{2,3\}} \int \left[ \frac{\delta b(\mathbf{x}_1) g(\mathbf{x}_1)}{\sqrt{\hat{g}(\mathbf{x}_1)}} \Pi \left( \frac{\delta a g}{\sqrt{\hat{g}}} | V_j^{\perp} \right)(\mathbf{x}_1) \right] d\mathbf{x}_1 \qquad (A.2)
$$

Plugging in A.2 into A.1, we get,

$$
|\mathrm{E}_P \left( \hat{\phi}_{n,j} - \phi(P) \right)|
$$

$$
= \left| \begin{array}{c} E_{P,\{2,3\}} \int \left[ \delta a(\mathbf{x}_1) \delta b(\mathbf{x}_1) g^2(\mathbf{x}_1) \left( \frac{1}{\hat{g}(\mathbf{x}_1)} - \frac{1}{g(\mathbf{x}_1)} \right) \right] d\mathbf{x}_1 \\ - E_{P,\{2,3\}} \int \left[ \frac{\delta b(\mathbf{x}_1) g(\mathbf{x}_1)}{\sqrt{\hat{g}(\mathbf{x}_1)}} \Pi \left( \frac{\delta a g}{\sqrt{\hat{g}}} | V_j^{\perp} \right)(\mathbf{x}_1) \right] d\mathbf{x}_1 \end{array} \right|
$$

$$(A.3)$$

Now, by repeatedly applying Cauchy-Schwarz Inequality and invoking results in Theorem 2.3, we have

$$
\left| E_{P,\{2,3\}} \int \left[ \delta a(\mathbf{x}_1) \delta b(\mathbf{x}_1) g^2(\mathbf{x}_1) \left( \frac{1}{\hat{g}(\mathbf{x}_1)} - \frac{1}{g(\mathbf{x}_1)} \right) \right] d\mathbf{x}_1 \right|
$$

$$
\leq \left( \mathrm{E}_{P,\{3\}} \int \frac{g^4(\mathbf{x}_1)}{g(\mathbf{x}_1) \hat{g}(\mathbf{x}_1)} (\hat{g}(\mathbf{x}_1) - g(\mathbf{x}_1))^2 d\mathbf{x}_1 \right)^{\frac{1}{2}}
$$

$$
\times \left( \mathrm{E}_{P,\{2,3\}} \int (\hat{a}(\mathbf{x}_1) - a(\mathbf{x}_1))^4 d\mathbf{x}_1 \right)^{\frac{1}{4}} \left( \mathrm{E}_{P,\{2,3\}} \int \left( \hat{b}(\mathbf{x}_1) - b(\mathbf{x}_1) \right)^4 d\mathbf{x}_1 \right)^{\frac{1}{4}}
$$

$$
\leq \frac{B_U^2}{B_L} \left( \mathrm{E}_{P,\{3\}} \|\hat{g} - g\|_2^2 \right)^{\frac{1}{2}} \left( \mathrm{E}_{P,\{2,3\}} \|\hat{a} - a\|_4^4 \right)^{\frac{1}{4}} \left( \mathrm{E}_{P,\{2,3\}} \left\| \hat{b} - b \right\|_4^4 \right)^{\frac{1}{4}}
$$

$$
\leq \frac{B_U^2}{B_L} (C)^{\frac{d}{2\gamma+d} + \frac{d}{2\alpha+d} + \frac{d}{2\beta+d}} \left( \frac{n}{\log n} \right)^{-\frac{\alpha}{2\alpha+d} - \frac{\beta}{2\beta+d} - \frac{\gamma}{2\gamma+d}} \qquad (A.4)
$$

Moreover,

$$\left| E_{P,\{2,3\}} \int \left[ \frac{\delta b(\mathbf{x}_1) g(\mathbf{x}_1)}{\sqrt{\hat{g}(\mathbf{x}_1)}} \Pi \left( \frac{\delta a g}{\sqrt{\hat{g}}} |V_j^\perp \right) (\mathbf{x}_1) \right] d\mathbf{x}_1 \right|$$

$$= \left| E_{P,\{2,3\}} \int \left[ \Pi \left( \frac{\delta b g}{\sqrt{\hat{g}}} |V_j^\perp \right) (\mathbf{x}_1) \Pi \left( \frac{\delta a g}{\sqrt{\hat{g}}} |V_j^\perp \right) (\mathbf{x}_1) \right] d\mathbf{x}_1 \right|$$

$$\leq \left( E_{P,\{2,3\}} \left[ \left\| \Pi \left( \frac{\delta b g}{\sqrt{\hat{g}}} |V_j^\perp \right) \right\|_2^2 \right] \right)^{\frac{1}{2}} \left( E_{P,\{2,3\}} \left[ \left\| \Pi \left( \frac{\delta a g}{\sqrt{\hat{g}}} |V_j^\perp \right) \right\|_2^2 \right] \right)^{\frac{1}{2}}$$

$$\leq C \left( 2^{-2j\beta} + \frac{1}{n^2} \right)^{\frac{1}{2}} \left( 2^{-2j\alpha} + \frac{1}{n^2} \right)^{\frac{1}{2}}, \tag{A.5}$$

where the last line follows for some constant $C$ (depending on $M, B_U, B_L, \gamma_{\max}$) by Theorem 2.3, definition of (5.5), and noting that $\| \Pi (h|V_j) \|_\infty \leq C(B_U)$ if $\|h\|_\infty \leq B_U$. Therefore, if $n \leq 2^{jd} \leq n^2$ along with $\frac{\alpha+\beta}{2} < \frac{d}{4}$, one has combining A.3, A.4, and A.5, that for a constant $C$ (depending on $M, B_U, B_L, \gamma_{\min}, \gamma_{\max}$) and $\gamma_{\min}(\epsilon) := \frac{\gamma_{\min}}{1+\epsilon}$

$$|\mathrm{E}_P \left( \hat{\phi}_{n,j} - \phi(P) \right)|$$

$$\leq C \left[ \left( \frac{n}{\log n} \right)^{-\frac{\alpha}{2\alpha+d} - \frac{\beta}{2\beta+d} - \frac{\gamma}{2\gamma+d}} + \left( 2^{-2j\beta} + \frac{1}{n^2} \right)^{\frac{1}{2}} \left( 2^{-2j\alpha} + \frac{1}{n^2} \right)^{\frac{1}{2}} \right]$$

$$\leq C \left[ \left( \frac{n}{\log n} \right)^{-\frac{\alpha}{2\alpha+d} - \frac{\beta}{2\beta+d} - \frac{\gamma}{2\gamma+d}} + 2^{-2jd\frac{\alpha+\beta}{2d}} + 3n^{-\frac{3}{2}} \right]$$

$$\leq 4C \left[ n^{-\frac{\alpha}{2\alpha+d} - \frac{\beta}{2\beta+d} - \frac{\gamma_{\min}(\epsilon)}{2\gamma_{\min}(\epsilon)+d}} n^{\frac{\gamma_{\min}(\epsilon)}{2\gamma_{\min}+d(\epsilon)} - \frac{\gamma}{2\gamma+d}} \log n^{\frac{\alpha}{2\alpha+d} + \frac{\beta}{2\beta+d} + \frac{\gamma}{2\gamma+d}} + 2^{-2jd\frac{\alpha+\beta}{2d}} \right]$$

$$\leq 4C \left[ n^{-\frac{\alpha}{2\alpha+d} - \frac{\beta}{2\beta+d} - \frac{\gamma_{\min}(\epsilon)}{2\gamma_{\min}(\epsilon)+d}} n^{-\frac{\gamma - \gamma_{\min}(\epsilon)}{2\gamma+d}} \log n + 2^{-2jd\frac{\alpha+\beta}{2d}} \right]$$

$$\leq 4C \left[ n^{-\frac{\alpha}{2\alpha+d} - \frac{\beta}{2\beta+d} - \frac{\gamma_{\min}(\epsilon)}{2\gamma_{\min}(\epsilon)+d}} n^{-\frac{\epsilon\gamma_{\min}}{(1+\epsilon)(2\max+d)}} \log n + 2^{-2jd\frac{\alpha+\beta}{2d}} \right]$$

Now, letting $\theta = (\alpha, \beta, \gamma)$, $f_1(\theta) = \frac{\alpha+\beta}{2}$ and $f_2(\theta) = -\frac{\alpha}{2\alpha+d} - \frac{\beta}{2\beta+d} - \frac{\gamma_{\min}(\epsilon)}{2\gamma_{\min}(\epsilon)+d}$ we have the bias property corresponding to Theorem 2.1 holds with the given choice of $f_1$ and $f_2$ and a constant $C$ depending on $M, B_U, B_L, \gamma_{\max}$ for $\{P \in \mathcal{P}_\theta : f_1(\theta) = \frac{\alpha+\beta}{2}, f_2(\theta) > \frac{2\alpha+2\beta}{2\alpha+2\beta+d}\}$. To proof of the validity of the variance property corresponding to Theorem 2.1 is easy to derive by standard Hoeffding decomposition of $\hat{\phi}_{n,j}$ followed by applications of moment bounds in Lemmas B.2 and B.5. For calculations of similar flavor, refer to proof of Theorem 1.3 in Mukherjee and Sen (2016). Note that this is the step where we have used the fact that $\frac{\alpha+\beta}{2} \leq \frac{d}{4}$, since otherwise the linear term dominates resulting in $O(\frac{1}{n})$ the variance.

$$\sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta) = \tau, f_2(\theta) > \frac{4\tau}{4\tau+d}}} \mathrm{E}_P \left( \hat{\phi}_{n,j(k^*(\hat{l}))} - \phi(P) \right)^2 \leq 8C \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{8\tau/d}{1+4\tau/d}}.$$

Noting that for $\theta \in \Theta$, since $\gamma_{\min} > 2(1+\epsilon) \max\{\alpha, \beta\}$, one has automatically, $f_2(\theta) > \frac{2\alpha+2\beta}{2\alpha+2\beta+d}$, completes the proof of the upper bound.

### (ii) Proof of Lower Bound

To prove a lower bound matching the upper bound above, note that $\phi(P) = \mathrm{E}_P\left(cov_P\left(Y, A|\mathbf{X}\right)\right) = \mathrm{E}_P\left(AY\right) - \mathrm{E}_P\left(a(\mathbf{X})b(\mathbf{X})\right)$. Indeed, $\mathrm{E}_P\left(AY\right)$ can be estimated at a $\sqrt{n}$-rate by sample average of $A_iY_i$. Therefore, it suffices to prove a lower for adaptive estimation of $\mathrm{E}_P\left(a(\mathbf{X})b(\mathbf{X})\right)$. Let $c(\mathbf{X}) = \mathrm{E}_P\left(Y|A=1, \mathbf{X}\right) - \mathrm{E}_P\left(Y|A=0, \mathbf{X}\right)$, which implies owing to the binary nature of $A$ that $\mathrm{E}_P\left(Y|A, \mathbf{X}\right) = c(\mathbf{X})\left(A - a(\mathbf{X})\right) + b(\mathbf{X})$. For the purpose of lower bound it is convenient to parametrize the data generating mechanism by $(a, b, c, g)$, which implies that $\phi(P) = \int a(\mathbf{x})b(\mathbf{x})g(\mathbf{x})d\mathbf{x}$. With this parametrization, we show that the same lower bound holds in a smaller class of problems where $g \equiv 1$ on $[0,1]^d$. Specifically consider

$$\Theta_{sub} = \left\{ \begin{array}{c} P = (a, b, c, g): \\ a \in H(\alpha, M), b \in H(\beta, M),\ \frac{\alpha+\beta}{2} < \frac{d}{4}, \\ g \equiv 1, (a(\mathbf{x}), b(\mathbf{x})) \in [B_L, B_U]^2\ \forall \mathbf{x} \in [0,1]^d \end{array} \right\}.$$

The likelihood of $O \sim P$ for $P \in \Theta_{sub}$ can then be written as

$$\begin{aligned} &a(\mathbf{X})^A(1 - a(\mathbf{X}))^{1-A} \\ &\times \left(c(\mathbf{X})(1 - A(X)) + b(X)\right)^{YA}\left(1 - c(\mathbf{X})(1 - a(\mathbf{X})) - b(\mathbf{X})\right)^{(1-Y)A} \\ &\times \left(-c(\mathbf{X})a(\mathbf{X}) + b(\mathbf{X})\right)^{Y(1-A)}\left(1 + c(\mathbf{X})a(\mathbf{X}) - b(\mathbf{X})\right)^{(1-Y)(1-A)}. \end{aligned} \tag{A.6}$$

Let for some $(\alpha, \beta, \gamma)$ tuple in the original problem $\Theta$, one has

$$\sup_{P \in \mathcal{P}_{(\alpha, \beta, \gamma)}} \mathrm{E}_P\left(\hat{\phi} - \phi(P)\right)^2 \leq C \left(\frac{\sqrt{\log n}}{n}\right)^{\frac{4\alpha+4\beta}{d+2\alpha+2\beta}}.$$

Now, let $H: [0,1]^d \to \mathbb{R}$ be a $C^\infty$ function supported on $\left[0, \frac{1}{2}\right]^d$ such that $\int H(\mathbf{x})d\mathbf{x} = 0$ and $\int H^2(\mathbf{x})d\mathbf{x} = 1$ and let for $k \in \mathbb{N}$ (to be decided later) $\Omega_1, \ldots, \Omega_k$ be the translates of the cube $k^{-\frac{1}{d}}\left[0, \frac{1}{2}\right]^d$ that are disjoint and contained in $[0,1]^d$. Let $\mathbf{x}_1, \ldots, \mathbf{x}_k$ denote the bottom left corners of these cubes.

Assume first that $\alpha < \beta$. We set for $\lambda = (\lambda_1, \ldots, \lambda_k) \in \{-1, +1\}^k$ and $\alpha \leq \beta' < \beta$,

$$a_\lambda(\mathbf{x}) = \frac{1}{2} + \left(\frac{1}{k}\right)^{\frac{\alpha}{d}}\sum_{j=1}^{k}\lambda_j H\left((\mathbf{x} - \mathbf{x}_j)k^{\frac{1}{d}}\right),$$

$$b_\lambda(\mathbf{x}) = \frac{1}{2} + \left(\frac{1}{k}\right)^{\frac{\beta'}{d}}\sum_{j=1}^{k}\lambda_j H\left((\mathbf{x} - \mathbf{x}_j)k^{\frac{1}{d}}\right),$$

$$c_\lambda(\mathbf{x}) = \frac{\frac{1}{2} - b_\lambda(\mathbf{x})}{1 - a_\lambda(\mathbf{x})}.$$

A properly chosen $H$ guarantees $a_\lambda \in H(\alpha, M)$ and $b_\lambda \in H(\beta', M)$ for all $\lambda$. Let

$$\Theta_0 = \left\{ P^n: P = \left(a_\lambda, \frac{1}{2}, 0, 1\right), \lambda \in \{-1, +1\}^k \right\},$$

and

$$\Theta_1 = \left\{ P^n: P = (a_\lambda, b_\lambda, c_\lambda, 1), \lambda \in \{-1, +1\}^k \right\}.$$

Finally let $\Theta_{test} = \Theta_0 \cup \Theta_1$. Let $\pi_0$ and $\pi_1$ be uniform priors on $\Theta_0$ and $\Theta_1$ respectively. It is easy to check that by our choice of $H$, $\phi(P) = \frac{1}{4}$ on $\Theta_0$ and $\phi(P) = \frac{1}{4} + \left(\frac{1}{k}\right)^{\frac{\alpha+\beta'}{d}}$ for $P \in \Theta_1$. Therefore, using notation from Lemma B.1, $\mu_1 = \frac{1}{4}$, $\mu_2 = \frac{1}{4} + \left(\frac{1}{k}\right)^{\frac{\alpha+\beta'}{d}}$, and $\sigma_1 = \sigma_2 = 0$. Since $\Theta_0 \subseteq P(\alpha, \beta, \gamma)$, we must have that worst case error of estimation over $\Theta_0$ is bounded by $C\left(\frac{\sqrt{\log n}}{n}\right)^{\frac{4\alpha+4\beta}{d+2\alpha+2\beta}}$. Therefore, the $\pi_0$ average bias over $\Theta_0$ is also bounded by $C\left(\frac{\sqrt{\log n}}{n}\right)^{\frac{2\alpha+2\beta}{d+2\alpha+2\beta}}$. This implies by Lemma B.1, that the $\pi_1$ average bias over $\Theta_1$ (and hence the worst case bias over $\Theta_1$) is bounded below by

$$\left(\frac{1}{k}\right)^{\frac{\alpha+\beta'}{d}} - C\left(\frac{\sqrt{\log n}}{n}\right)^{\frac{2\alpha+2\beta}{d+2\alpha+2\beta}} - C\left(\frac{\sqrt{\log n}}{n}\right)^{\frac{2\alpha+2\beta}{d+2\alpha+2\beta}}\eta, \tag{A.7}$$

where $\eta$ is the chi-square divergence between the probability measures $\int P^n d\pi_0(P^n)$ and $\int P^n d\pi_1(P^n)$. We now bound $\eta$ using Theorem 2.2.

To put ourselves in the notation of Theorem 2.2, let for $\lambda \in \{-1, +1\}^k$, $P_\lambda$ and $Q_\lambda$ be the probability measures identified from $\Theta_0$ and $\Theta_1$ respectively.

Therefore, with $\chi_j = \{0, 1\} \times \{0, 1\} \times \Omega_j$, we indeed have for all $j = 1, \dots, k$, $P_\lambda(\chi_j) = Q_\lambda(\chi_j) = p_j$ where there exists a constant $c$ such that $p_j = \frac{c}{k}$.

Letting $\pi$ be the uniform prior over $\{-1, +1\}^k$ it is immediate that $\eta = \chi^2\left(\int P_\lambda d(\pi(\lambda)), \int Q_\lambda d(\pi(\lambda))\right)$.

It now follows by calculations similar to proof of Theorem 4.1 in Robins et al. (2009), that for a constant $C' > 0$

$$\chi^2\left(\int P_\lambda d(\pi(\lambda)), \int Q_\lambda d(\pi(\lambda))\right) \leq \exp\left(C'\frac{n^2}{k}\left(k^{-\frac{4\beta'}{d}} + k^{-4\frac{\alpha+\beta'}{2d}}\right)\right) - 1.$$

Now choosing $k = \left(\frac{n}{\sqrt{c_* \log n}}\right)^{\frac{2d}{d+2\alpha+2\beta'}}$, we have

$$\chi^2\left(\int P_\lambda d(\pi(\lambda)), \int Q_\lambda d(\pi(\lambda))\right) \leq n^{2C'c_*} - 1.$$

Therefore choosing $c_*$ such that $2C'c_* + \frac{2\alpha+2\beta'}{2\alpha+2\beta'+d} < \frac{2\alpha+2\beta}{2\alpha+2\beta+d}$, we have the desired result by (A.7). The proof for $\alpha > \beta$ is similar after changing various quantities to:

$$a_\lambda(\mathbf{x}) = \frac{1}{2} + \left(\frac{1}{k}\right)^{\frac{\alpha'}{d}}\sum_{j=1}^{k}\lambda_j H\left((\mathbf{x} - \mathbf{x}_j)k^{\frac{1}{d}}\right), \quad \beta \leq \alpha' < \alpha,$$

$$b_\lambda(\mathbf{x}) = \frac{1}{2} + \left(\frac{1}{k}\right)^{\frac{\beta}{d}}\sum_{j=1}^{k}\lambda_j H\left((\mathbf{x} - \mathbf{x}_j)k^{\frac{1}{d}}\right),$$

$$c_\lambda(\mathbf{X}) = \frac{(\frac{1}{2} - a_\lambda(\mathbf{X}))b_\lambda(\mathbf{X})}{a_\lambda(\mathbf{X})(1 - a_\lambda(\mathbf{X}))},$$

$$\Theta_0 = \left\{P^n \colon P = \left(\frac{1}{2}, b_\lambda, 0, 1\right) : \lambda \in \{-1, +1\}^k\right\},$$

and

$$\Theta_1 = \left\{ P^n \colon P = (a_\lambda, b_\lambda, c_\lambda, 1) \colon \lambda \in \{-1, +1\}^k \right\}.$$

For the case of $\alpha = \beta$, choose $\alpha' < \beta$ and therefore, $\alpha' < \beta$ and thereafter work with

$$a_\lambda(\mathbf{x}) = \frac{1}{2} + \left(\frac{1}{k}\right)^{\frac{\alpha'}{d}} \sum_{j=1}^{k} \lambda_j H\left((\mathbf{x} - \mathbf{x}_j)k^{\frac{1}{d}}\right),$$

$$b_\lambda(\mathbf{x}) = \frac{1}{2} + \left(\frac{1}{k}\right)^{\frac{\beta}{d}} \sum_{j=1}^{k} \lambda_j H\left((\mathbf{x} - \mathbf{x}_j)k^{\frac{1}{d}}\right),$$

$$c_\lambda(\mathbf{x}) = \frac{\frac{1}{2} - b_\lambda(\mathbf{x})}{1 - a_\lambda(\mathbf{x})}.$$

$$\Theta_0 = \left\{ P^n \colon P = \left(a_\lambda, \frac{1}{2}, 0, 1\right) \colon \lambda \in \{-1, +1\}^k \right\},$$

and

$$\Theta_1 = \left\{ P^n \colon P = (a_\lambda, b_\lambda, c_\lambda, 1) \colon \lambda \in \{-1, +1\}^k \right\}.$$

This completes the proof of the lower bound. ∎

**Proof of Theorem 3.2.**

PROOF. *(i) Proof of Upper Bound*

The general scheme of proof is same as that of Theorem 3.2 and involves identifying a non-adaptive minimax estimator of $\phi(P)$ under the knowledge of $P \in \mathcal{P}_{(\alpha,\beta,\gamma)}$, demonstrating suitable bias and variance properties of this sequence of estimators, and thereafter invoking Theorem 2.1 to conclude. This routine can be carried out as follows. Without loss of generality assume that we have $3n$ samples $\{Y_i A_i, A_i, \mathbf{X}_i\}_{i=1}^{n}$. Divide the samples in to 3 equal parts (with the $l^{\text{th}}$ part being indexed by $\{(l-1)n+1, \ldots, ln\}$ for $l \in \{1, 2, 3\}$), estimate $f$ by $\hat{f}$ adaptively from the third part (as in Theorem 2.3), and estimate $\mathrm{E}(A|\mathbf{x})$ and $b$ by $\widehat{\mathrm{E}(A|\mathbf{x})}$ and $\hat{b}(\mathbf{x}) := \widehat{\mathrm{E}(Y|A=1, \mathbf{x})}$ respectively, adaptively from the second part (as in Theorem 2.3). Let $\mathrm{E}_{P,S}$ denote the expectation while samples with indices in $S$ held fixed, for $S \subset \{1, 2, 3\}$. Note that $g(\mathbf{X}) = f(\mathbf{X}|A=1)P(A=1)$. Therefore, also estimate $P(A=1)$ by $\hat{\pi} := \frac{1}{n}\sum_{2n+1}^{3n} A_i$ i.e. the sample average of $A$'s from the third part of the sample and $\hat{f}_1$ is estimated as an estimator of $f(X|A=1)$ from the third part of our sample using density estimation technique among observations with $A = 1$. Finally, our estimate of $a$ and $g$ are $\hat{a}(\mathbf{x}) = \frac{1}{\widehat{\mathrm{E}(A|\mathbf{x})}}$ and $\hat{g} = \hat{f}_1\hat{\pi}$ respectively. In the following, we will freely use Theorem 2.3, for desired properties of $\hat{a}, \hat{b}$, and $\hat{g}$. In particular, following the proof of Theorem 2.3, we can actually assume that our choice of $\hat{g}$ also satisfies the necessary conditions of boundedness away from 0 and $\infty$, as well as membership in $H(\gamma, C)$ with high probability for a large enough $C > 0$. A first order influence function for $\phi(P)$ at $P$ is given by $Aa(\mathbf{X})(Y - b(\mathbf{X}) + b(\mathbf{X}) - \phi(P)$ and a resulting first order estimator for $\phi(P)$ is $\frac{1}{n}\sum_{i=1}^{n} A_i a(\mathbf{X}_i)(Y_i - \hat{b}(\mathbf{X}_i)) + b(\mathbf{X}_i)$. This estimator has a bias $-\mathrm{E}_{P,\{2,3\}} \int \left((b(\mathbf{x}) - \hat{b}(\mathbf{x}))(a(\mathbf{x}) - \hat{a}(\mathbf{x}))\right) g(\mathbf{x})d\mathbf{x}$. Indeed for $\frac{\alpha+\beta}{2} < \frac{d}{2}$, this bias turns out to

be suboptimal compared to the minimax rate of convergence of $n^{-\frac{4\alpha+4\beta}{2\alpha+2\beta+d}}$ in mean squared loss. Similar to proof of Theorem 3.1 we use a second order bias corrected estimator as follows.

Once again we work with the linear subspace $V_j$ (defined in 5) where the choice of $j$ is guided by the balance between the bias and variance properties of the resulting estimator. In particular, a choice of $2^j$ is guided by the knowledge of the parameter space $\mathcal{P}_{(\alpha,\beta,\gamma)}$. For any $j$ such that $n \leq 2^{jd} \leq n^2$, our bias corrected second order estimator of $\phi(P)$ is given by

$$\hat{\phi}_{n,j} = \frac{1}{n} \sum_{i=1}^{n} A_i \hat{a}(\mathbf{X}_i)(Y_i - \hat{b}(\mathbf{X}_i)) + \hat{b}(\mathbf{X}_i)$$
$$+ \frac{1}{n(n-1)} \sum_{1 \leq i_1 \neq i_2 \leq n} S\left( \frac{A_{i_1}(Y_{i_1} - \hat{b}(\mathbf{X}_{i_1})))}{\sqrt{\hat{g}(\mathbf{X}_{i_1})}} K_{V_j}(\mathbf{X}_{i_1}, \mathbf{X}_{i_2}) \frac{(A_{i_2}\hat{a}(\mathbf{X}_{i_2}) - 1)}{\sqrt{\hat{g}(\mathbf{X}_{i_2})}} \right)$$

Note that division by $\hat{g}$ is permitted by the properties guaranteed by Theorem 2.3. Indeed this sequence of estimators is in the form of those considered by Theorem 2.1 with

$$L_1(O) = A\hat{a}(\mathbf{X})(Y - \hat{b}(\mathbf{X})) + \hat{b}(\mathbf{X}),$$

$$L_{2l}(O) = -\frac{A(Y - \hat{b}(\mathbf{X})))}{\sqrt{\hat{g}(\mathbf{X})}}, \quad L_{2r}(O) = \frac{(A\hat{a}(\mathbf{X}) - 1)}{\sqrt{\hat{g}(\mathbf{X})}},$$

where by Theorem 2.3 $\max\{|L_1(O)|, |L_{2l}(O)|, |L_{2r}(O)|\} \leq C(B_L, B_U)$. Therefore it remains to show that these sequence $\hat{\phi}_{n,j}$ satisfies the bias and variance property (A) and (B) necessary for application of Theorem 2.1. Using the conditional independence of $Y$ and $A$ given $\mathbf{X}$, one has y calculations exactly parallel to that in proof of Theorem 3.1, that for a constant $C$ (depending on $M, B_U, B_L, \gamma_{\min}, \gamma_{\max}$),

$$|\mathrm{E}_P\left(\hat{\phi}_{n,j} - \phi(P)\right)|$$
$$\leq C\left[ n^{-\frac{\alpha}{2\alpha+d} - \frac{\beta}{2\beta+d} - \frac{\gamma_{\min}(\epsilon)}{2\gamma_{\min}(\epsilon)+d}} n^{-\frac{\epsilon\gamma_{\min}}{(1+\epsilon)(2\max+d)}} \log n + 2^{-2jd\frac{\alpha+\beta}{2d}} \right],$$

where $\gamma_{\min}(\epsilon) := \frac{\gamma_{\min}}{1+\epsilon}$. Now, letting $\theta = (\alpha, \beta, \gamma)$, $f_1(\theta) = \frac{\alpha+\beta}{2}$ and $f_2(\theta) = -\frac{\alpha}{2\alpha+d} - \frac{\beta}{2\beta+d} - \frac{\gamma_{\min}(\epsilon)}{2\gamma_{\min}(\epsilon)+d}$ we have the bias property corresponding to Theorem 2.1 holds with the given choice of $f_1$ and $f_2$ and a constant $C$ depending on $M, B_U, B_L, \gamma_{\max}$ for $\{P \in \mathcal{P}_\theta : f_1(\theta) = \frac{\alpha+\beta}{2}, f_2(\theta) > \frac{2\alpha+2\beta}{2\alpha+2\beta+d}\}$. To proof of the validity of the variance property corresponding to Theorem 2.1 is one again easy to derive by standard Hoeffding decomposition of $\hat{\phi}_{n,j}$ followed by applications of moment bounds in Lemmas B.2 and B.5.

$$\sup_{\substack{P \in \mathcal{P}_\theta: \\ f_1(\theta)=\tau, f_2(\theta) > \frac{4\tau}{4\tau+d}}} \mathrm{E}_P\left(\hat{\phi}_{n,j\left(k^*(\hat{l})\right)} - \phi(P)\right)^2 \leq 8C\left(\frac{\sqrt{\log n}}{n}\right)^{\frac{8\tau/d}{1+4\tau/d}}.$$

Noting that for $\theta \in \Theta$, since $\gamma_{\min} > 2(1+\epsilon)\max\{\alpha, \beta\}$, one has automatically, $f_2(\theta) > \frac{2\alpha+2\beta}{2\alpha+2\beta+d}$, completes the proof of the upper bound.

### (ii) Proof of Lower Bound
First note that we can parametrize our distributions by the tuple of functions $(a, b, g)$. We show

that the same lower bound holds in a smaller class of problems where $g \equiv 1/2$ on $[0,1]^d$. Specifically consider

$$\Theta_{sub} = \left\{ \begin{array}{c} P = (a,b,g): \\ a \in H(\alpha, M), b \in H(\beta, M), \frac{\alpha+\beta}{2} < \frac{d}{4}, \\ g \equiv 1/2, (a(\mathbf{x}), b(\mathbf{x})) \in [B_L, B_U]^2 \ \forall \mathbf{x} \in [0,1]^d \end{array} \right\}.$$

The observed data likelihood of $O \sim P$ for $P \in \Theta_{sub}$ can then be written as

$$(a(\mathbf{X}) - 1)^{1-A} \left( b^Y(\mathbf{X})(1 - b(\mathbf{X}))^{1-Y} \right)^A. \tag{A.8}$$

Let for some $(\alpha, \beta, \gamma)$ tuple in the original problem $\Theta$, one has

$$\sup_{P \in \mathcal{P}_{(\alpha,\beta,\gamma)}} \mathrm{E}_P \left( \hat{\phi} - \phi(P) \right)^2 \leq C \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{4\alpha+4\beta}{d+2\alpha+2\beta}}.$$

Now, let $H: [0,1]^d \to \mathbb{R}$ be a $C^\infty$ function supported on $\left[0, \frac{1}{2}\right]^d$ such that $\int H(\mathbf{x})d\mathbf{x} = 0$ and $\int H^2(\mathbf{x})d\mathbf{x} = 1$ and let for $k \in \mathbb{N}$ (to be decided later) $\Omega_1, \ldots, \Omega_k$ be the translates of the cube $k^{-\frac{1}{d}} \left[0, \frac{1}{2}\right]^d$ that are disjoint and contained in $[0,1]^d$. Let $\mathbf{x}_1, \ldots, \mathbf{x}_k$ denote the bottom left corners of these cubes.

Assume first that $\alpha < \beta$. We set for $\lambda = (\lambda_1, \ldots, \lambda_k) \in \{-1, +1\}^k$ and $\alpha \leq \beta' < \beta$,

$$a_\lambda(\mathbf{x}) = 2 + \left( \frac{1}{k} \right)^{\frac{\alpha}{d}} \sum_{j=1}^k \lambda_j H\left( (\mathbf{x} - \mathbf{x}_j)k^{\frac{1}{d}} \right),$$

$$b_\lambda(\mathbf{x}) = \frac{1}{2} + \left( \frac{1}{k} \right)^{\frac{\beta'}{d}} \sum_{j=1}^k \lambda_j H\left( (\mathbf{x} - \mathbf{x}_j)k^{\frac{1}{d}} \right).$$

A properly chosen $H$ guarantees $a_\lambda \in H(\alpha, M)$ and $b_\lambda \in H(\beta', M)$ for all $\lambda$. Let

$$\Theta_0 = \left\{ P^n: P = (a_\lambda, 1/2, 1/2) : \lambda \in \{-1, +1\}^k \right\},$$

and

$$\Theta_1 = \left\{ P^n: P = (a_\lambda, b_\lambda, 1/2): \lambda \in \{-1, +1\}^k \right\}.$$

Finally let $\Theta_{test} = \Theta_0 \cup \Theta_1$. Let $\pi_0$ and $\pi_1$ be uniform priors on $\Theta_0$ and $\Theta_1$ respectively. It is easy to check that by our choice of $H$, $\phi(P) = \frac{1}{2}$ on $\Theta_0$ and $\phi(P) = \frac{1}{2} + \frac{1}{2} \left( \frac{1}{k} \right)^{\frac{\alpha+\beta'}{d}}$ for $P \in \Theta_1$. Therefore, using notation from Lemma B.1, $\mu_1 = \frac{1}{2}$, $\mu_2 = \frac{1}{2} + \frac{1}{2} \left( \frac{1}{k} \right)^{\frac{\alpha+\beta'}{d}}$, and $\sigma_1 = \sigma_2 = 0$. Since $\Theta_0 \subseteq P(\alpha, \beta, \gamma)$, we must have that worst case error of estimation over $\Theta_0$ is bounded by $C \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{4\alpha+4\beta}{d+2\alpha+2\beta}}$. Therefore, the $\pi_0$ average bias over $\Theta_0$ is also bounded by $C \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{2\alpha+2\beta}{d+2\alpha+2\beta}}$. This implies by Lemma B.1, that the $\pi_1$ average bias over $\Theta_1$ (and hence the worst case bias over $\Theta_1$) is bounded below by a constant multiple of

$$\left( \frac{1}{k} \right)^{\frac{\alpha+\beta'}{d}} - \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{2\alpha+2\beta}{d+2\alpha+2\beta}} - \left( \frac{\sqrt{\log n}}{n} \right)^{\frac{2\alpha+2\beta}{d+2\alpha+2\beta}} \eta, \tag{A.9}$$

where $\eta$ is the chi-square divergence between the probability measures $\int P^n d\pi_0(P^n)$ and $\int P^n d\pi_1(P^n)$. We now bound $\eta$ using Theorem 2.2.

To put ourselves in the notation of Theorem 2.2, let for $\lambda \in \{-1, +1\}^k$, $P_\lambda$ and $Q_\lambda$ be the probability measures identified from $\Theta_0$ and $\Theta_1$ respectively.

Therefore, with $\chi_j = \{0, 1\} \times \{0, 1\} \times \Omega_j$, we indeed have for all $j = 1, \ldots, k$, $P_\lambda(\chi_j) = Q_\lambda(\chi_j) = p_j$ where there exists a constant $c$ such that $p_j = \frac{c}{k}$.

Letting $\pi$ be the uniform prior over $\{-1, +1\}^k$ it is immediate that $\eta = \chi^2 \left( \int P_\lambda d(\pi(\lambda)), \int Q_\lambda d(\pi(\lambda)) \right)$.

It now follows by calculations similar to proof of Theorem 4.1 in Robins et al. (2009), that for a constant $C' > 0$

$$\chi^2 \left( \int P_\lambda d(\pi(\lambda)), \int Q_\lambda d(\pi(\lambda)) \right) \leq \exp\left( C' \frac{n^2}{k} \left( k^{-\frac{4\beta'}{d}} + k^{-4\frac{\alpha+\beta'}{2d}} \right) \right) - 1.$$

Now choosing $k = \left( \frac{n}{\sqrt{c_* \log n}} \right)^{\frac{2d}{d+2\alpha+2\beta'}}$, we have

$$\chi^2 \left( \int P_\lambda d(\pi(\lambda)), \int Q_\lambda d(\pi(\lambda)) \right) \leq n^{2C'c_*} - 1.$$

Therefore choosing $c_*$ such that $2C'c_* + \frac{2\alpha+2\beta'}{2\alpha+2\beta'+d} < \frac{2\alpha+2\beta}{2\alpha+2\beta+d}$, we have the desired result by (A.7). The proof for $\alpha > \beta$ is similar after changing various quantities to:

$$a_\lambda(\mathbf{x}) = 2 + \left( \frac{1}{k} \right)^{\frac{\alpha'}{d}} \sum_{j=1}^{k} \lambda_j H\left( (\mathbf{x} - \mathbf{x}_j) k^{\frac{1}{d}} \right), \quad \beta \leq \alpha' < \alpha,$$

$$b_\lambda(\mathbf{x}) = \frac{1}{2} + \left( \frac{1}{k} \right)^{\frac{\beta}{d}} \sum_{j=1}^{k} \lambda_j H\left( (\mathbf{x} - \mathbf{x}_j) k^{\frac{1}{d}} \right).$$

$$\Theta_0 = \left\{ P^n \colon P = (2, b_\lambda, 1/2) : \lambda \in \{-1, +1\}^k \right\},$$

and

$$\Theta_1 = \left\{ P^n \colon P = (a_\lambda, b_\lambda, 1/2) : \lambda \in \{-1, +1\}^k \right\}.$$

For the case of $\alpha = \beta$, choose $\alpha' < \beta$ and therefore, $\alpha' < \beta$ and thereafter work with

$$a_\lambda(\mathbf{x}) = 2 + \left( \frac{1}{k} \right)^{\frac{\alpha'}{d}} \sum_{j=1}^{k} \lambda_j H\left( (\mathbf{x} - \mathbf{x}_j) k^{\frac{1}{d}} \right),$$

$$b_\lambda(\mathbf{x}) = \frac{1}{2} + \left( \frac{1}{k} \right)^{\frac{\beta}{d}} \sum_{j=1}^{k} \lambda_j H\left( (\mathbf{x} - \mathbf{x}_j) k^{\frac{1}{d}} \right)$$

$$\Theta_0 = \left\{ P^n \colon P = (a_\lambda, 1/2, 1/2) : \lambda \in \{-1, +1\}^k \right\},$$

and

$$\Theta_1 = \left\{ P^n \colon P = (a_\lambda, b_\lambda, 1/2) : \lambda \in \{-1, +1\}^k \right\}.$$

This completes the proof of the lower bound.

∎

**Proof of Theorem 3.3.**

PROOF. *(i) Proof of Upper Bound*

Without loss of generality assume that we have $3n$ samples $\{Y_i, A_i, \mathbf{X}_i\}_{i=1}^n$. Divide the samples in to 3 equal parts (with the $l^{\text{th}}$ part being indexed by $\{(l-1)n+1, \ldots, ln\}$ for $l \in \{1, 2, 3\}$), estimate $g$ by $\hat{g}$ adaptively from the third part (as in Theorem 2.3), and estimate $b$ by $\hat{b}$, adaptively from the second part (as in Theorem 2.3). Let $\mathrm{E}_{P,S}$ denote the expectation while samples with indices in $S$ held fixed, for $S \subset \{1, 2, 3\}$. For any $j$ such that $n \le 2^{jd} \le n^2$, consider

$$\hat{\phi}_{n,j} = \frac{1}{n} \sum_{i=1}^n (2Y_i - \hat{b}(\mathbf{X}_i))\hat{b}(\mathbf{X}_i)$$
$$+ \frac{1}{n(n-1)} \sum_{1 \le i_1 \ne i_2 \le n} \frac{(Y_{i_1} - \hat{b}(\mathbf{X}_{i_1})))}{\sqrt{\hat{g}(\mathbf{X}_{i_1})}} K_{V_j}(\mathbf{X}_{i_1}, \mathbf{X}_{i_2}) \frac{(Y_{i_2} - \hat{b}(\mathbf{X}_{i_2}))}{\sqrt{\hat{g}(\mathbf{X}_{i_2})}}$$

Indeed this sequence of estimators is in the form of those considered by Theorem 2.1 with

$$L_1(O) = (2Y - \hat{b}(\mathbf{X}))\hat{b}(\mathbf{X}),$$

$$L_{2l}(O) = -\frac{(Y - \hat{b}(\mathbf{X})))}{\sqrt{\hat{g}(\mathbf{X})}}, \quad L_{2r}(O) = \frac{(Y - \hat{b}(\mathbf{X})))}{\sqrt{\hat{g}(\mathbf{X})}},$$

where by Theorem 2.3 $\max\{|L_1(O)|, |L_{2l}(O)|, |L_{2r}(O)|\} \le C(B_L, B_U)$. Therefore it remains to show that these sequence $\hat{\phi}_{n,j}$ satisfies the bias and variance property (A) and (B) necessary for application of Theorem 2.1. We first verify the bias property. Utilizing the representation of the first order bias as stated above, we have

$$|\mathrm{E}_P\left(\hat{\phi}_{n,j} - \phi(P)\right)|$$
$$= \left| \begin{array}{c} \mathrm{E}_{P,\{2,3\}}\left[\int \left((b(\mathbf{x}) - \hat{b}(\mathbf{x})\right)^2 g(\mathbf{x})d\mathbf{x}\right] \\ -\mathrm{E}_P\left[S\left(\frac{(Y_1 - \hat{b}(\mathbf{X}_1)))}{\sqrt{\hat{g}(\mathbf{X}_1)}} K_{V_j}(\mathbf{X}_1, \mathbf{X}_2) \frac{(Y_2 - \hat{b}(\mathbf{X}_2))}{\sqrt{\hat{g}(\mathbf{X}_2)}}\right)\right] \end{array} \right|$$

(A.10)

Now, by calculations similar to proof of Theorem 3.1, one can show that for a constant $C$ (depending on $M, B_U, B_L, \gamma_{\min}, \gamma_{\max}$),

$$|\mathrm{E}_P\left(\hat{\phi}_{n,j} - \phi(P)\right)|$$
$$\le C\left[n^{-\frac{2\beta}{2\beta+d} - \frac{\gamma_{\min}(\epsilon)}{2\gamma_{\min}(\epsilon)+d}} n^{-\frac{\epsilon\gamma_{\min}}{(1+\epsilon)(2\max+d)}} \log n + 2^{-2jd\frac{\alpha+\beta}{2d}}\right], \quad \gamma_{\min}(\epsilon) := \frac{\gamma_{\min}}{1+\epsilon}$$

Now, letting $\theta = (\beta, \gamma)$, $f_1(\theta) = \beta$ and $f_2(\theta) = -\frac{2\beta}{2\beta+d} - \frac{\gamma_{\min}(\epsilon)}{2\gamma_{\min}(\epsilon)+d}$, the rest of the proof follows along the lines of the proof of Theorem 3.1.

*(ii) Proof of Lower Bound*

The proof of the lower bound is very similar to that of the lower bound proof in Theorem 3.1, after identifying $Y = A$ almost surely $\mathbb{P}$, and hence is omitted. ∎

## APPENDIX B: TECHNICAL LEMMAS

**B.1. Constrained Risk Inequality.** A main tool for producing adaptive lower bound arguments is a general version of constrained risk inequality due to Cai and Low (2011), obtained as an extension of Brown and Low (1996). For the sake of completeness, begin with a summary of these results. Suppose $Z$ has distribution $\mathbb{P}_\theta$ where $\theta$ belongs to some parameter space $\Theta$. Let $\hat{Q} = \hat{Q}(Z)$ be an estimator of a function $Q(\theta)$ based on $Z$ with bias $B(\theta) := \mathrm{E}_\theta(\hat{Q}) - Q(\theta)$. Now suppose that $\Theta_0$ and $\Theta_1$ form a disjoint partition of $\Theta$ with priors $\pi_0$ and $\pi_1$ supported on them respectively. Also, let $\mu_i = \int Q(\theta) d\pi_i$ and $\sigma_i^2 = \int (Q(\theta) - \mu_i)^2 d\pi_i$, $i = 0, 1$ be the mean and variance of $Q(\theta)$ under the two priors $\pi_0$ and $\pi_1$. Letting $\gamma_i$ be the marginal density with respect to some common dominating measure of $Z$ under $\pi_i$, $i = 0, 1$, let us denote by $\mathrm{E}_{\gamma_0}(g(Z))$ the expectation of $g(Z)$ with respect to the marginal density of $Z$ under prior $\pi_0$ and distinguish it from $\mathrm{E}_\theta(g(Z))$, which is the expectation under $\mathbb{P}_\theta$. Lastly, denote the chi-square divergence between $\gamma_0$ and $\gamma_1$ by $\chi = \left\{ \mathrm{E}_{\gamma_0} \left( \frac{\gamma_1}{\gamma_0} - 1 \right)^2 \right\}^{\frac{1}{2}}$. Then we have the following result.

LEMMA B.1 (Cai and Low (2011)). *If* $\int \mathrm{E}_\theta \left( \hat{Q}(Z) - Q(\theta) \right)^2 d\pi_0(\theta) \leq \epsilon^2$, *then*

$$\left| \int B(\theta) d\pi_1(\theta) - \int B(\theta) d\pi_0(\theta) \right| \geq |\mu_1 - \mu_0| - (\epsilon + \sigma_0)\chi.$$

Since the maximum risk is always at least as large as the average risk, this immediately yields a lower bound on the minimax risk.

**B.2. Tail and Moment Bounds.** The U-statistics appearing in this paper are mostly based on projection kernels sandwiched between arbitrary bounded functions. This necessitates generalizing the U-statistics bounds obtained in Bull and Nickl (2013) as in Mukherjee and Sen (2016).

LEMMA B.2. $\mathbf{O}_1, \ldots, \mathbf{O}_n \sim \mathbb{P}$ *are iid random vectors of observations such that* $\mathbf{X}_i \in [0,1]^d$ *is a sub-vector of* $\mathbf{O}_i$ *for each* $i$. *There exists constant* $C := C(B, B_U, J_0) > 0$ *such that the following hold*

*(i)*

$$(\text{B.1}) \qquad \mathbb{P}\left( \left| \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} R(\mathbf{O}_{i_1}, \mathbf{O}_{i_2}) - \mathrm{E}(R(\mathbf{O}_1, \mathbf{O}_2)) \right| \geq t \right)$$

$$\leq e^{-Cnt^2} + e^{-\frac{Ct^2}{a_1^2}} + e^{-\frac{Ct}{a_2}} + e^{-\frac{C\sqrt{t}}{\sqrt{a_3}}},$$

*(ii)*

$$\mathrm{E}\left( \left| \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} R(\mathbf{O}_{i_1}, \mathbf{O}_{i_2}) - \mathrm{E}(R(\mathbf{O}_1, \mathbf{O}_2)) \right|^{2q} \right) \leq \left( C \frac{2^{jd}}{n^2} \right)^q,$$

*where* $a_1 = \frac{1}{n-1} 2^{\frac{jd}{2}}$, $a_2 = \frac{1}{n-1} \left( \sqrt{\frac{2^{jd}}{n}} + 1 \right)$, $a_3 = \frac{1}{n-1} \left( \sqrt{\frac{2^{jd}}{n}} + \frac{2^{jd}}{n} \right)$,
$R(\mathbf{O}_1, \mathbf{O}_2) = S \left( L_{2l}(O_1) K_{V_j}(\mathbf{X}_1, \mathbf{X}_2) L_{2r}(O_2) \right)$ *with* $\max\{|L_{2l}(O)|, |L_{2r}(O)|\} \leq B$, *almost surely* $\mathbf{O}$, *and* $\mathbf{X}_i \in [0,1]^d$ *are iid with density* $g$ *such that* $g(\mathbf{x}) \leq B_U$ *for all* $\mathbf{x} \in [0,1]^d$.

PROOF. The proof of part (i) can be found in Mukherjee and Sen (2016). However, for the sake of completeness we provide the proof here again. We do the proof for the special case where $L_{2l} = L_{2r} = L$. However, the details of the argument shows that the proof continues to hold to symmetrized U-statistics as defined here.

The proof hinges on the following tail bound for second order degenerate U-statistics (Giné and Nickl, 2015) is due to Giné, Latala and Zinn (2000) with constants by Houdré and Reynaud-Bouret (2003) and is crucial for our calculations.

LEMMA B.3. *Let $U_n$ be a degenerate U-statistic of order 2 with kernel $R$ based on an i.i.d. sample $W_1, \ldots, W_n$. Then there exists a constant $C$ independent of $n$, such that*

$$P[|\sum_{i \neq j} R(W_1, W_2)| \geq C(\Lambda_1 \sqrt{u} + \Lambda_2 u + \Lambda_3 u^{3/2} + \Lambda_4 u^2)] \leq 6 \exp(-u),$$

*where, we have,*

$$\Lambda_1^2 = \frac{n(n-1)}{2} E[R^2(W_1, W_2)],$$

$$\Lambda_2 = n \sup\{E[R(W_1, W_2)\zeta(W_1)\xi(W_2)] : E[\zeta^2(W_1)] \leq 1, E[\xi^2(W_1)] \leq 1\},$$

$$\Lambda_3 = \|nE[R^2(W_1, \cdot)]\|_\infty^{\frac{1}{2}},$$

$$\Lambda_4 = \|R\|_\infty.$$

We use this lemma to establish Lemma B.2. By Hoeffding's decomposition one has

$$\frac{1}{n(n-1)} \sum_{i_1 \neq i_2} R(\mathbf{O}_{i_1}, \mathbf{O}_{i_2}) - \mathrm{E}(R(\mathbf{O}_1, \mathbf{O}_2))$$

$$= \frac{2}{n} \sum_{i_1=1}^{n} \left[ \mathrm{E}_{\mathbf{O}_{i_1}} R(\mathbf{O}_{i_1}, \mathbf{O}_{i_2}) - \mathrm{E}R(\mathbf{O}_{i_1}, \mathbf{O}_{i_2}) \right]$$

$$+ \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} \left[ \begin{array}{c} R(\mathbf{O}_{i_1}, \mathbf{O}_{i_2}) - \mathrm{E}_{\mathbf{O}_{i_1}} R(\mathbf{O}_{i_1}, \mathbf{O}_{i_2}) \\ -\mathrm{E}_{\mathbf{O}_{i_2}} R(\mathbf{O}_{i_1}, \mathbf{O}_{i_2}) + \mathrm{E}R(\mathbf{O}_{i_1}, \mathbf{O}_{i_2}) \end{array} \right]$$

$$:= T_1 + T_2$$

B.2.1. *Analysis of $T_1$.* Noting that $T_1 = \frac{2}{n} \sum_{i_1=1}^{n} H(\mathbf{O}_{i_1})$ where $H(\mathbf{O}_{i_1}) = \mathrm{E}(R(\mathbf{O}_{i_1}, \mathbf{O}_{i_2}|\mathbf{O}_{i_1})) - \mathrm{E}R(\mathbf{O}_{i_1}, \mathbf{O}_{i_2})$ we control $T_1$ by standard Hoeffding's Inequality. First note that,

$$|H(\mathbf{O}_{i_1})|$$

$$= |\sum_{k \in \mathcal{Z}_j} \sum_{v \in \{0,1\}^d} \left[ L(\mathbf{O}_{i_1}) \psi_{jk}^v(\mathbf{X}_{i_1}) \mathrm{E}(\psi_{jk}^v(\mathbf{X}_{i_2}) L(\mathbf{O}_{i_2})) - (\mathrm{E}(\psi_{jk}^v(\mathbf{X}_{i_2}) L(\mathbf{O}_{i_2})))^2 \right]|$$

$$\leq \sum_{k \in \mathcal{Z}_j} \sum_{v \in \{0,1\}^d} |L(\mathbf{O}_{i_1}) \psi_{jk}^v(\mathbf{X}_{i_1}) \mathrm{E}(\psi_{jk}^v(\mathbf{X}_{i_2}) L(\mathbf{O}_{i_2}))|$$

$$+ \sum_{k \in \mathcal{Z}_j} \sum_{v \in \{0,1\}^d} (\mathrm{E}(\psi_{jk}^v(\mathbf{X}_{i_2}) L(\mathbf{O}_{i_2})))^2$$

First, by standard compactness argument for the wavelet bases,

$$|\mathrm{E}(\psi_{jk}^v(\mathbf{X}) L(\mathbf{O}))| \leq \int |\mathrm{E}(L(\mathbf{O})|\mathbf{X} = \mathbf{x}) \left( 2^{\frac{jd}{2}} \prod_{l=1}^{d} \psi_{00}^{v_l}(2^j x_l - k_l) \right)||g(\mathbf{x})|d\mathbf{x}$$

$$\leq C(B, B_U, J_0) 2^{-\frac{jd}{2}}. \tag{B.2}$$

Therefore,

$$\sum_{k \in \mathcal{Z}_j} \sum_{v \in \{0,1\}^d} \left( \mathrm{E}\left( \psi_{jk}^v\left( \mathbf{X}_{i_2} \right) L\left( \mathbf{O}_{i_2} \right) \right) \right)^2 \leq C(B, B_U, J_0) \tag{B.3}$$

Also, using the fact that for each fixed $\mathbf{x} \in [0,1]^d$, the number indices $k \in \mathcal{Z}_j$ such that $\mathbf{x}$ belongs to support of at least one of $\psi_{jk}^v$ is bounded by a constant depending only on $\psi_{00}^0$ and $\psi_{00}^1$. Therefore combining (B.2) and (B.3),

$$\sum_{k \in \mathcal{Z}_j} \sum_{v \in \{0,1\}^d} \left| L\left( \mathbf{O}_{i_1} \right) \psi_{jk}^v\left( \mathbf{X}_{i_1} \right) \mathrm{E}\left( \psi_{jk}^v\left( \mathbf{X}_{i_2} \right) L\left( \mathbf{O}_{i_2} \right) \right) \right|$$

$$\leq C(B, B_U, J_0) 2^{-\frac{jd}{2}} 2^{\frac{jd}{2}} = C(B, B_U, J_0). \tag{B.4}$$

Therefore, by (B.4) and Hoeffding's Inequality,

$$\mathbb{P}\left( |T_1| \geq t \right) \leq 2 e^{-C(B, B_U, J_0) n t^2}. \tag{B.5}$$

B.2.2. *Analysis of $T_2$.* Since $T_2$ is a degenerate U-statistics, it's analysis is based on Lemma B.3. In particular,

$$T_2 = \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} R^*\left( \mathbf{O}_{i_1}, \mathbf{O}_{i_2} \right)$$

where

$$R^*\left( \mathbf{O}_{i_1}, \mathbf{O}_{i_2} \right)$$

$$= \sum_{k \in \mathcal{Z}_j} \sum_{v \in \{0,1\}^d} \left\{ \begin{array}{l} \left( L(\mathbf{O}_{i_1}) \psi_{jk}^v\left( \mathbf{X}_{i_1} \right) - \mathrm{E}\left( \psi_{jk}^v\left( \mathbf{X}_{i_1} \right) \mathrm{E}\left( L(\mathbf{O}_{i_1}) | \mathbf{X}_{i_1} \right) \right) \right) \\ \times \left( L(\mathbf{O}_{i_2}) \psi_{jk}^v\left( \mathbf{X}_{i_2} \right) - \mathrm{E}\left( \psi_{jk}^v\left( \mathbf{X}_{i_2} \right) \mathrm{E}\left( L(\mathbf{O}_{i_2}) | \mathbf{X}_{i_2} \right) \right) \right) \end{array} \right\}$$

Letting $\Lambda_i$, $i = 1, \ldots, 4$ being the relevant quantities as in Lemma B.3, we have the following lemma.

LEMMA B.4. *There exists a constant $C = C(B, B_U, J_0)$ such that*

$$\Lambda_1^2 \leq C \frac{n(n-1)}{2} 2^{jd}, \ \Lambda_2 \leq Cn, \ \Lambda_3^2 \leq Cn 2^{jd}, \ \Lambda_4 \leq C 2^{\frac{jd}{2}}.$$

PROOF. First we control $\Lambda_1$. To this end, note that by simple calculations, using bounds on $L, g$, and orthonormality of $\psi_{jk}^v$'s we have,

$$\Lambda_1^2 = \frac{n(n-1)}{2} \mathrm{E}\left( \{ R^*\left( \mathbf{O}_1, \mathbf{O}_2 \right) \}^2 \right) \leq 3n(n-1) \mathrm{E}\left( R^2\left( \mathbf{O}_1, \mathbf{O}_2 \right) \right)$$

$$= 3n(n-1) \mathrm{E}\left( L^2\left( \mathbf{O}_1 \right) K_{V_j}^2\left( \mathbf{X}_1, \mathbf{X}_2 \right) L^2\left( \mathbf{O}_2 \right) \right)$$

$$\leq 3n(n-1) B^4 \int \int \left[ \sum_{k \in \mathcal{Z}_j} \sum_{v \in \{0,1\}^d} \psi_{jk}^v\left( \mathbf{x}_1 \right) \psi_{jk}^v\left( \mathbf{x}_2 \right) \right]^2 g(\mathbf{x}_1) g(\mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2$$

$$\leq 3n(n-1)B^4 B_U^2 \int\int \Big[\sum_{k\in\mathcal{Z}_j}\sum_{v\in\{0,1\}^d} \psi_{jk}^v(\mathbf{x}_1)\,\psi_{jk}^v(\mathbf{x}_2)\Big]^2 d\mathbf{x}_1 d\mathbf{x}_2$$

$$= 3n(n-1)B^4 B_U^2 \sum_{k\in\mathcal{Z}_j}\sum_{v\in\{0,1\}^d}\int \big(\psi_{jk}^v(\mathbf{x}_1)\big)^2 d\mathbf{x}_2 \int \big(\psi_{jk}^v(\mathbf{x}_2)\big)^2 d\mathbf{x}_2$$

$$\leq C(B,B_U,J_0)n(n-1)2^{jd}.$$

Next we control

$$\Lambda_2 = n\sup\big\{\mathrm{E}\left(R^*(\mathbf{O}_1,\mathbf{O}_2)\zeta(\mathbf{O}_1)\xi(\mathbf{O}_2)\right): \mathrm{E}\left(\zeta^2(\mathbf{O}_1)\right)\leq 1, \mathrm{E}\left(\xi^2(\mathbf{O}_2)\right)\leq 1\big\}.$$

To this end, we first control

$$|\mathrm{E}\left(L(\mathbf{O}_1)K_{V_j}(\mathbf{X}_1,\mathbf{X}_2)L(\mathbf{O}_2)\zeta(\mathbf{O}_1)\xi(\mathbf{O}_2)\right)|$$
$$= |\int\int \mathrm{E}(L(\mathbf{O}_1)\zeta(\mathbf{O}_1)|\mathbf{X}_1=\mathbf{x}_1)K_{V_j}(\mathbf{x}_1,\mathbf{x}_2)\,\mathrm{E}(L(\mathbf{O}_2)\xi(\mathbf{O}_2)|\mathbf{X}_2=\mathbf{x}_2)g(\mathbf{x}_2)g(\mathbf{x}_2)d\mathbf{x}_1 d\mathbf{x}_2|$$
$$= |\int \mathrm{E}(L(\mathbf{O})\zeta(\mathbf{O})|\mathbf{X}=\mathbf{x})\Pi\left(\mathrm{E}(L(\mathbf{O})\xi(\mathbf{O})|\mathbf{X}=\mathbf{x})g(\mathbf{x})|V_j\right)g(\mathbf{x})d\mathbf{x}|$$
$$\leq \left(\int \mathrm{E}^2(L(\mathbf{O})\zeta(\mathbf{O})|\mathbf{X}=\mathbf{x})g^2(\mathbf{x})d\mathbf{x}\right)^{\frac12}\left(\int \Pi^2\left(\mathrm{E}(L(\mathbf{O})\xi(\mathbf{O})|\mathbf{X}=\mathbf{x})g(\mathbf{x})|V_j\right)d\mathbf{x}\right)^{\frac12}$$
$$\leq \left(\int \mathrm{E}(L^2(\mathbf{O})\zeta^2(\mathbf{O})|\mathbf{X}=\mathbf{x})g^2(\mathbf{x})d\mathbf{x}\right)^{\frac12}\left(\int \mathrm{E}(L^2(\mathbf{O})\xi^2(\mathbf{O})|\mathbf{X}=\mathbf{x})g^2(\mathbf{x})d\mathbf{x}\right)^{\frac12}$$
$$\leq B^2 B_U\sqrt{\mathrm{E}(\zeta^2(\mathbf{O}_1))\mathrm{E}(\xi^2(\mathbf{O}_2))}\leq B^2 B_U$$

Above we have used Cauchy-Schwartz Inequality, Jensen's Inequality, and the fact that projections contract norm. Also,

$$|\mathrm{E}\left(\mathrm{E}\left(L(\mathbf{O}_1)K_{V_j}(\mathbf{X}_1,\mathbf{X}_2)L(\mathbf{O}_2)|\mathbf{O}_1\right)\zeta(\mathbf{O}_1)\xi(\mathbf{O}_2)\right)|$$
$$= |\mathrm{E}\left[L(\mathbf{O}_1)\Pi\left(\mathrm{E}\left(L(\mathbf{O}_1)g(\mathbf{X}_1)|\mathbf{X}_1\right)|V_j\right)\zeta(\mathbf{O}_1)\xi(\mathbf{O}_2)\right]|$$
$$= |\mathrm{E}\left[L(\mathbf{O}_1)\Pi\left(\mathrm{E}\left(L(\mathbf{O}_1)g(\mathbf{X}_1)|\mathbf{X}_1\right)|V_j\right)\zeta(\mathbf{O}_1)\right]||\mathrm{E}(\xi(\mathbf{O}_2))|$$
$$\leq |\int \Pi(\mathrm{E}(L(\mathbf{O})\zeta(\mathbf{O})|\mathbf{X}=\mathbf{x})g(\mathbf{x})|V_j)\Pi(\mathrm{E}(L(\mathbf{O})|\mathbf{X}=\mathbf{x})g(\mathbf{x})|V_j)d\mathbf{x}|\leq B^2 B_U,$$

where the last step once again uses contraction property of projection, Jensen's Inequality, and bounds on $L$ and $g$. Finally, by Cauchy-Schwartz Inequality and (B.3),

$$\mathrm{E}\left[\mathrm{E}\left(L(\mathbf{O}_1)K_{V_j}(\mathbf{X}_1,\mathbf{X}_2)L(\mathbf{O}_2)\right)\zeta(\mathbf{O}_1)\xi(\mathbf{O}_2)\right]$$
$$\leq \sum_{k\in\mathcal{Z}_j}\sum_{v\in\{0,1\}^d}\mathrm{E}^2\left(L(\mathbf{O})\psi_{jk}^v(\mathbf{X})\right)\leq C(B,B_U,J_0).$$

This completes the proof of $\Lambda_2\leq C(B,B_U,J_0)n$. Turning to $\Lambda_3 = n\|\mathrm{E}\left[(R^*(\mathbf{O}_1,\cdot))^2\right]\|_\infty^{\frac12}$ we have that

$$(R^*(\mathbf{O}_1,\mathbf{o}_2))^2$$
$$\leq 2\left[R(\mathbf{O}_1,\mathbf{o}_2)-\mathrm{E}(R(\mathbf{O}_1,\mathbf{O}_2)|\mathbf{O}_1)\right]^2+2\left[\mathrm{E}(R(\mathbf{O}_1,\mathbf{O}_2)|\mathbf{O}_2=\mathbf{o}_2)-\mathrm{E}\left(R(\mathbf{O}_1,\mathbf{O}_2)\right)\right]^2$$

Now,

$$\mathrm{E}\left[R(\mathbf{O}_1,\mathbf{o}_2)-\mathrm{E}(R(\mathbf{O}_1,\mathbf{O}_2)|\mathbf{O}_1)\right]^2$$

$$\leq 2\mathrm{E}\left(L^2(\mathbf{O}_1)K_{V_j}^2\left(\mathbf{X}_1,\mathbf{x}_2\right)L^2(\mathbf{o}_2)\right)$$

$$+ 2\mathrm{E}\Big(\sum_{k\in\mathcal{Z}_j}\sum_{v\in\{0,1\}^d}L(\mathbf{O}_1)\psi_{jk}^v(\mathbf{X}_1)\mathrm{E}\left(\psi_{jk}^v(\mathbf{X}_2)L(\mathbf{O}_2)\right)\Big)^2$$

$$\leq 2B^4B_U^2\sum_{k\in\mathcal{Z}_j}\sum_{v\in\{0,1\}^d}\left(\psi_{jk}^v(\mathbf{x}_2)\right)^2 + 2\mathrm{E}(H^2(\mathbf{O}_2)) \leq C(B,B_U,J_0)2^{jd}.$$

where the last inequality follows from arguments along the line of (B.4). Also, using inequalities (B.3) and (B.4)

$$[\mathrm{E}(R(\mathbf{O}_1,\mathbf{O}_2)|\mathbf{O}_2=\mathbf{o}_2) - \mathrm{E}\left(R(\mathbf{O}_1,\mathbf{O}_2)\right)]^2$$

$$= \Big[\sum_{k\in\mathcal{Z}_j}\sum_{v\in\{0,1\}^d}\mathrm{E}\left(L(\mathbf{O}_1)\psi_{jk}^v(\mathbf{X}_1)\right)\left(\mathrm{E}\left(L(\mathbf{O}_1)\psi_{jk}^v(\mathbf{X}_1)\right) - \psi_{jk}^v(\mathbf{x}_2)L(\mathbf{o}_2)\right)\Big]^2$$

$$\leq C(B,B_U,J_0).$$

This completes the proof of controlling $\Lambda_3$. Finally, using compactness of the wavelet basis,

$$\|R(\cdot,\cdot)\|_\infty \leq B^2\sup_{\mathbf{x}_1,\mathbf{x}_2}\sum_{k\in\mathcal{Z}_j}\sum_{v\in\{0,1\}^d}|\psi_{jk}^v(\mathbf{x}_1)||\psi_{jk}^v(\mathbf{x}_2)| \leq C(B,B_U,J_0)2^{jd}$$

Combining this with arguments similar to those leading to (B.4), we have $\Lambda_4 \leq C(B,B_U,J_0)2^{jd}$. ∎

Therefore, using Lemma B.3 and Lemma B.4 we have

$$\mathbb{P}\Big(|T_2| \geq \frac{C(B,B_U,J_0)}{n-1}\Big(\sqrt{2^{jd}t} + t + \sqrt{\frac{2^{jd}}{n}}t^{\frac{3}{2}} + \frac{2^{jd}}{n}t^2\Big)\Big) \leq 6e^{-t}.$$

Finally using $2t^{\frac{3}{2}} \leq t + t^2$ we have,

$$\mathbb{P}_f\left[|T_2| > a_1\sqrt{t} + a_2t + a_3t^2\right] \leq 6e^{-t} \tag{B.6}$$

where $a_1 = \frac{C(B,B_U,J_0)}{n-1}2^{\frac{jd}{2}}$, $a_2 = \frac{C(B,B_U,J_0)}{n-1}\left(\sqrt{\frac{2^{jd}}{n}} + 1\right)$,

$a_3 = \frac{C(B,B_U,J_0)}{n-1}\left(\sqrt{\frac{2^{jd}}{n}} + \frac{2^{jd}}{n}\right)$. Now if $h(t)$ is such that $a_1\sqrt{h(t)} + a_2h(t) + a_3h^2(t) \leq t$, then one has by (B.6),

$$\mathbb{P}\left[|T_2| \geq t\right] \leq \mathbb{P}\left[|T_2| \geq a_1\sqrt{h(t)} + a_2h(t) + a_3h^2(t)\right] \leq 6e^{-6h(t)}.$$

Indeed, there exists such an $h(t)$ such that $h(t) = b_1t^2 \wedge b_2t \wedge b_3\sqrt{t}$ where $b_1 = \frac{C(B,B_U,J_0)}{a_1^2}$, $b_2 = \frac{C(B,B_U,J_0)}{a_2}$, and $b_3 = \frac{C(B,B_U,J_0)}{\sqrt{a_3}}$. Therefore, there exists $C = C(B,B_U,J_0)$ such that

$$\mathbb{P}\left[|T_2| \geq t\right] \leq e^{-\frac{Ct^2}{a_1^2}} + e^{-\frac{Ct}{a_2}} + e^{-\frac{C\sqrt{t}}{\sqrt{a_3}}}. \tag{B.7}$$

B.2.3. *Combining Bounds on $T_1$ and $T_2$.* Applying union bound along with B.5 and B.7 completes the proof of Lemma B.2 part (i).

For the proof of part (ii) note that with the notation of the proof of part (i) we have by Hoeffding decomposition

$$\mathrm{E}\left(\Big|\frac{1}{n(n-1)}\sum_{i_1\neq i_2}R\left(\mathbf{O}_{i_1},\mathbf{O}_{i_2}\right) - \mathrm{E}\left(R\left(\mathbf{O}_1,\mathbf{O}_2\right)\right)\Big|^{2q}\right) \leq 2(\mathrm{E}|T_1|^{2q} + \mathrm{E}|T_2|^{2q})$$

The proof will be completed by individual control of the two moments above.

B.2.4. *Analysis of* $\mathrm{E}|T_1|^{2q}$. Recall that $T_1 = \frac{2}{n}\sum_{i_1=1}^{n} H(\mathbf{O}_{i_1})$ where $H(\mathbf{O}_{i_1}) = \mathrm{E}\left(R\left(\mathbf{O}_{i_1}, \mathbf{O}_{i_2}|\mathbf{O}_{i_1}\right)\right) - \mathrm{E}R\left(\mathbf{O}_{i_1}, \mathbf{O}_{i_2}\right)$ and $|H(\mathbf{O})| \leq C(B, B_U, J_0)$ almost surely. Therefore by Rosenthal's Inequality B.5 we have

$$\mathrm{E}|T_1|^{2q} \leq \left(\frac{2}{n}\right)^{2q}\left[\sum_{i=1}^{n}\mathrm{E}|H(\mathbf{O}_i)|^{2q} + \left\{\sum_{i=1}^{n}\mathrm{E}|H(\mathbf{O}_i)|^2\right\}^q\right]$$

$$\leq \left(\frac{2C(B, B_U, J_0)}{n}\right)^{2q}(n + n^q) \leq C(B, B_U, J_0)^q n^{-q}.$$

B.2.5. *Analysis of* $\mathrm{E}|T_2|^{2q}$. Recall that

$$\mathbb{P}\left[|T_2| \geq t\right] \leq \mathbb{P}\left[|T_2| \geq a_1\sqrt{h(t)} + a_2 h(t) + a_3 h^2(t)\right] \leq 6e^{-6h(t)}.$$

where $h(t) = b_1 t^2 \wedge b_2 t \wedge b_3\sqrt{t}$ with $b_1 = \frac{C(B, B_U, J_0)}{a_1^2}$, $b_2 = \frac{C(B, B_U, J_0)}{a_2}$, and $b_3 = \frac{C(B, B_U, J_0)}{\sqrt{a_3}}$. Therefore

$$\mathbb{E}_f(|T_2|^{2q})$$

$$= 2q\int_0^\infty x^{2q-1}\mathbb{P}_f(|T_2| \geq x)dx$$

$$\leq 2q\int_0^\infty x^{2q-1}\mathbb{P}_f(|T_2| \geq a_1\sqrt{h(x)} + a_2 h(x) + a_3 h^2(x))dx$$

$$\leq 12q\int_0^\infty x^{2q-1}e^{-h(x)}dx$$

$$= 12q\int_0^\infty x^{2q-1}e^{-\left\{b_1 x^2 \wedge b_2 x \wedge b_3\sqrt{x}\right\}}dx$$

$$\leq 12q\left[\int_0^\infty x^{2q-1}e^{-b_1 x^2}dx + \int_0^\infty x^{2q-1}e^{-b_2 x}dx + \int_0^\infty x^{2q-1}e^{-b_3\sqrt{x}}dx\right]$$

$$= 12q\left(\frac{\Gamma(q)}{2b_1^q} + \frac{\Gamma(2q)}{b_2^{2q}} + \frac{2\Gamma(4q)}{b_3^{4q}}\right) \leq \left(C\frac{2^{jd}}{n^2}\right)^q$$

for a constant $C = C(B, B_U, J_0)$, by our choices of $b_1, b_2, b_3$. ∎

Since the estimators arising in this paper also have a linear term, we will need the following standard Bernstein and Rosenthal type tail and moment bounds (Petrov, 1995).

LEMMA B.5. *If* $\mathbf{O}_1, \ldots, \mathbf{O}_n \sim \mathbb{P}$ *are iid random vectors such that* $|L(\mathbf{O})| \leq B$ *almost surely* $\mathbb{P}$, *then for* $q \geq 2$ *one has for large enough constants* $C(B)$ *and* $C(B, q)$

$$\mathbb{P}(|\frac{1}{n}\sum_{i=1}^{n}\left(L(\mathbf{O}_i) - \mathrm{E}(L(\mathbf{O}_i))\right)| \geq t) \leq 2e^{-nt^2/C(B)},$$

*and*

$$\mathrm{E}(|\sum_{i=1}^{n}\left(L(\mathbf{O}_i) - \mathrm{E}(L(\mathbf{O}_i))\right)|^q)$$

$$\leq \left[\sum_{i=1}^{n}\mathrm{E}\left(|L(\mathbf{O}_i) - \mathrm{E}(L(\mathbf{O}_i))|^q\right) + \left[\sum_{i=1}^{n}\mathrm{E}\left(|L(\mathbf{O}_i) - \mathrm{E}(L(\mathbf{O}_i))|^2\right)\right]^{q/2}\right]$$

$$\leq C(B, q)n^{\frac{q}{2}}.$$

We will also need the following concentration inequality for linear estimators based on wavelet projection kernels, proof of which can be done along the lines of proofs of Theorem 5.1.5 and Theorem 5.1.13 of Giné and Nickl (2015).

LEMMA B.6. *Consider i.i.d. observations* $\mathbf{O}_i = (Y, \mathbf{X})_i$, $i = 1, \ldots, n$ *where* $\mathbf{X}_i \in [0,1]^d$ *with marginal density* $g$. *Let* $\hat{m}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^{n} L(\mathbf{O}_i) K_{V_l}(\mathbf{X}_i, \mathbf{x})$, *such that* $\max\{\|g\|_\infty, \|L\|_\infty\} \leq B_U$. *If* $\frac{2^{ld}ld}{n} \leq 1$, *there exists* $C, C_1, C_2 > 0$, *depending on* $B_U$ *and scaling functions* $\psi_{0,0}^0, \psi_{0,0}^1$ *respectively, such that*

$$\mathrm{E}(\|\hat{m} - \mathrm{E}(\hat{m})\|_\infty) \leq C\sqrt{\frac{2^{ld}ld}{n}},$$

*and for any* $x > 0$

$$\mathbb{P}\left(n\|\hat{m} - \mathrm{E}(\hat{m})\|_\infty > \frac{3}{2}n\mathrm{E}(\|\hat{m} - \mathrm{E}(\hat{m})\|_\infty) + \sqrt{C_1 n 2^{ld} x} + C_2 2^{ld} x\right) \leq e^{-x}.$$