

# QUASI-BAYESIAN ANALYSIS OF NONPARAMETRIC INSTRUMENTAL VARIABLES MODELS\*

BY KENGO KATO

*University of Tokyo*

This paper aims at developing a quasi-Bayesian analysis of the nonparametric instrumental variables model, with a focus on the asymptotic properties of quasi-posterior distributions. In this paper, instead of assuming a distributional assumption on the data generating process, we consider a quasi-likelihood induced from the conditional moment restriction, and put priors on the function-valued parameter. We call the resulting posterior quasi-posterior, which corresponds to “Gibbs posterior” in the literature. Here we focus on priors constructed on slowly growing finite dimensional sieves. We derive rates of contraction and a non-parametric Bernstein-von Mises type result for the quasi-posterior distribution, and rates of convergence for the quasi-Bayes estimator defined by the posterior expectation. We show that, with priors suitably chosen, the quasi-posterior distribution (the quasi-Bayes estimator) attains the minimax optimal rate of contraction (convergence, respectively). These results greatly sharpen the previous related work.

## 1. Introduction.

1.1. *Overview.* Let  $(Y, X, W)$  be a triplet of scalar random variables, where  $Y$  is a dependent variable,  $X$  is an endogenous variable and  $W$  is an instrumental variable. Without losing much generality, we assume that the support of  $(X, W)$  is contained in  $[0, 1]^2$ . The support of  $Y$  may be unbounded. We consider the nonparametric instrumental variables (NPIV) model of the form

$$(1) \quad \mathbb{E}[Y | W] = \mathbb{E}[g_0(X) | W],$$

where  $g_0 : [0, 1] \rightarrow \mathbb{R}$  is an unknown structural function of interest. Alternatively we can write the model in a more conventional form

$$Y = g_0(X) + U, \quad \mathbb{E}[U | W] = 0,$$

---

\*Supported by the Grant-in-Aid for Young Scientists (B) (25780152) from the JSPS.  
*AMS 2000 subject classifications:* 62G08, 62G20

*Keywords and phrases:* asymptotic normality, inverse problem, nonparametric instrumental variables model, quasi-Bayes, rates of contraction

where  $X$  is potentially correlated with  $U$  and hence  $\mathbb{E}[U | X] \neq 0$ .

A model of the form (1) is of principal importance in econometrics [see 28, 33]. From a statistical perspective, the problem of recovering the structural function  $g_0$  is challenging since it is an *ill-posed* inverse problem with an additional difficulty of *unknown* operator ( $K$  in (2) ahead). Statistical inverse problems, including the current problem, have attracted considerable interests in statistics and econometrics [see, e.g., 9, 10]. For mathematical background of inverse problems, we refer to [45].

To see that the problem of recovering the structural function  $g_0$  is an ill-posed inverse problem, suppose that  $(X, W)$  has a square-integrable joint density  $f_{X,W}(x, w)$  on  $[0, 1]^2$  and denote by  $f_W(w)$  the marginal density of  $W$ . Define the linear operator  $K : L_2[0, 1] \rightarrow L_2[0, 1]$  by

$$(Kg)(w) = \mathbb{E}[g(X) | W = w]f_W(w) = \int g(x)f_{X,W}(x, w)dx.$$

Then the NPIV model (1) is equivalent to the operator equation

$$(2) \quad Kg_0 = h,$$

where  $h(w) = \mathbb{E}[Y | W = w]f_W(w)$ . Suppose that  $K$  is injective to guarantee identification of  $g_0$ .<sup>1</sup> The problem is that, even though  $K$  is injective, its inverse  $K^{-1}$  is not  $L_2$ -continuous since  $K$  is Hilbert-Schmidt (as  $f_{X,W}(x, w)$  is square integrable on  $[0, 1]^2$ ) and hence the  $l$ -th largest singular value, denoted by  $\kappa_l$ , is approaching zero as  $l \rightarrow \infty$  [see, e.g., 61]. In this sense, the problem of recovering  $g_0$  from  $h$  is ill-posed.

Approaches to estimating the structural function  $g_0$  are roughly classified into two types: the method involving the Tikhonov regularization [28, 17] and the sieve-based method [50, 2, 6, 34].<sup>2</sup> The minimax optimal rates of convergence in estimating  $g_0$  are established in [28, 12], and they are achieved by the estimators proposed in [28, 6] under their respective assumptions. All the above mentioned studies are, however, from a purely frequentist perspective. Little is known about the theoretical properties of Bayes or quasi-Bayes analysis of the NPIV model. Exceptions are [19, 20, 18, 46].

This paper aims at developing a quasi-Bayesian analysis of the NPIV model, with a focus on the asymptotic properties of quasi-posterior distributions. The approach taken is quasi-Bayes in the sense that it neither needs

<sup>1</sup>This global identification condition is, however, not a trivial assumption; see the discussion after Assumption 2 in Section 3.2 as well as the last paragraph in the next subsection.

<sup>2</sup>The sieve-method is further classified into two types: the method using slowly growing finite dimensional sieves with no or light penalties where the dimensions of sieves play the role of regularization, and the method using large dimensional sieves with heavy penalties where the penalty terms play the role of regularization [see 11].

to assume any specific distribution of  $(Y, X, W)$ , nor has to put a nonparametric prior on the unknown likelihood function. The analysis is then based upon a quasi-likelihood induced from the conditional moment restriction. The quasi-likelihood is constructed by first estimating the conditional moment function  $m(\cdot, g) = \mathbb{E}[Y - g(X) \mid W = \cdot]$  in a nonparametric way, and taking  $\exp\{-(1/2) \sum_{i=1}^n \hat{m}^2(W_i, g)\}$  as if it were a likelihood of  $g$ . For this quasi-likelihood, we put a prior on the function-valued parameter  $g$ . By doing so, formally, the posterior distribution for  $g$  may be defined, which we call “quasi-posterior distribution”. This posterior corresponds to what [37] called “Gibbs posterior”, and has a substantial interpretation (see Proposition 1 ahead). The quasi-Bayesian approach in this paper builds upon [13] where the dimension of the parameter of interest is finite and fixed.

We focus here on priors constructed on slowly growing finite dimensional sieves (called “sieve or series priors”), where the dimensions of the sieve spaces (which grow with the sample size) play the role of regularization to deal with the problem of ill-posedness. Potentially, there are several choices in sieve spaces, but we choose to use wavelet bases to form sieve spaces. Wavelet bases are useful to treat smoothness function classes such as Hölder-Zygmund and Sobolev spaces in a unified and convenient way. We also use wavelet series estimation of the conditional moment function.<sup>3</sup>

Under this setup, we study the asymptotic properties of the quasi-posterior distribution. The results obtained are summarized as follows. First, we derive rates of contraction for the quasi-posterior distribution and establish conditions on priors under which the minimax optimal rate of contraction is attained. Here the contraction is stated in the standard  $L_2$ -norm. Second, we show asymptotic normality of the quasi-posterior of the first  $k_n$  generalized Fourier coefficients, where  $k_n \rightarrow \infty$  is the dimension of the sieve space. This may be viewed as a non-parametric Bernstein-von Mises type result [see 59, Chapter 10 for the classical Bernstein-von Mises theorem for regular parametric models]. Third, we derive rates of convergence of the quasi-Bayes estimator defined by the posterior expectation and show that under some conditions it attains the minimax optimal rate of convergence. Finally, we give some specific sieve priors for which the quasi-posterior distribution (the quasi-Bayes estimator) attains the minimax optimal rate of contraction (convergence, respectively). These results greatly sharpen the previous work of e.g. [46], as we will review below.

---

<sup>3</sup>This does not rule out the use of other bases such as the Fourier and Hermite polynomial bases. See Remark 3.

1.2. *Literature review and contributions.* Closely related are [20] and [46]. The former paper worked on the reduced form equation  $Y = \mathbb{E}[g_0(X) | W] + V$  with  $V = U + g_0(X) - \mathbb{E}[g_0(X) | W]$  and assumed  $V$  to be normally distributed. They considered a Gaussian prior on  $g$ , and the posterior distribution is also Gaussian (conditionally on the variance of  $V$ ). They proposed to “regularize” the posterior and studied the asymptotic properties of the “regularized” posterior distribution and its expectation. Clearly, the present paper largely differs from [20] in that (i) we do not assume normality of the “error”; (ii) roughly speaking, Florens and Simoni’s method is tied with the Tikhonov regularization method, while ours is tied with the sieve-based method with slowly growing sieves. We note the settings of [19, 18] are largely different from the present paper; moreover in the NPIV example, some high-level conditions on estimated operators are assumed in [19, 18], and hence they are not directly comparable to the present paper. [46] developed an important unified framework in estimating conditional moment restriction models based on a quasi-Bayesian approach, and their scope is more general than ours. They analyzed NPIV models in detail in their Section 4. Their posterior construction is similar to ours such as the use of sieve priors, but differs from ours in detail. For example, [46] transformed the conditional moment restriction into unconditional moment restrictions with increasing number of restrictions. On the other hand, we directly work on the conditional moment restriction, although whether Liao and Jiang’s approach will lose any efficiency in the frequentist sense is not formally clear.

Importantly and substantially, neither [20] nor [46] established sharp contraction rates for their (quasi-)posterior distributions, nor asymptotic normality results. It is unclear whether Florens and Simoni’s [20] rates (in their Theorem 2) are optimal, since their assumptions are substantially different from the past literature such as [28] and [12]; moreover strictly speaking [20] did not formally derive contraction rates for their regularized posterior when the operator is unknown (note that [19, 18], though not directly comparable to the present paper, also did not formally derive posterior contraction rates in the NPIV example). [46] only established posterior consistency. Here we focus on a simple but important model, and establish the sharper asymptotic results for the quasi-posterior distribution. Notably, a wide class of (finite dimensional) sieve priors is shown to lead to the optimal contraction rate. Moreover, in [46], a point estimator of the structural function is not formally analyzed. Hence the primal contribution of this paper is to considerably deepen the understanding of the asymptotic properties of the quasi-Bayesian procedure for the NPIV model.

The present paper deals with a quasi-Bayesian analysis of an infinite di-

mensional model. The literature on theoretical studies of Bayesian analysis of infinite dimensional models is large. See [24, 56, 26, 40, 25] for general contraction rates results for posterior distributions in infinite dimensional models. Note that these results do not directly apply to our case: the proof of the main general theorem (Theorem 1) depends on the construction of suitable “tests” (see the proof of Proposition 4), but how to construct such tests in a specific problem in a non-likelihood framework is not trivial, especially in the current NPIV model where we have to deal with the ill-posedness of inverse problem. Moreover, Proposition 4 alone is not sufficient for obtaining sharp contraction rates and an additional work is needed (see the proof of Theorem 1).

There is also a large literature on the Bayesian analysis of (ill-posed) inverse problems. One stream of research on this topic lies in the applied mathematics literature; see [57] and references therein. However, their models and scopes are substantially different from those of the present paper; e.g., [31, 32] considered (ill-conditioned) finite-dimensional linear regression models with Gaussian errors and priors, and contractions rates of posterior distributions are not formally studied there. In the statistics literature, we may refer to [16, 42, 43, 1, 41] (in addition to [46, 19, 20, 18] that are already discussed), although their results are not applicable to the analysis of NPIV models because of its particular structure (i.e., especially the operator  $K$  is unknown, and non-Gaussian “errors” and priors are allowed). Hence the present paper provides a further contribution to the Bayesian analysis of ill-posed inverse problems.

Our asymptotic normality result builds upon the previous work on asymptotic normality of (quasi-)posterior distributions for models with increasing number of parameters [22, 23, 3, 4, 8, 14, 7]. Related is [7], in which the author established Bernstein-von Mises theorems for Gaussian regression models with increasing number of regressors and improved upon the earlier work of [22] in several aspects. [7] covered nonparametric models by taking into account modeling bias in the analysis. However, none of these papers covered the NPIV model, nor more generally linear inverse problems.

Lastly, while we here assume injectivity of the operator  $K$  in (2), as one of anonymous referees pointed out, this condition is not a trivial assumption (see also the discussion after Assumption 2 in Section 3.2), and there are a number of works that relax the injectivity assumption and explore partial identification approach, such as [52, 46, 44], and [11, Appendix A].

1.3. *Organization and notation.* The remainder of the paper is organized as follows. Section 2 gives an informal discussion of the quasi-Bayesian anal-

ysis of the NPIV model. Section 3 contains the main results of the paper where general theorems on contraction rates and asymptotic normality for quasi-posterior distributions, as well as convergence rates for quasi-Bayes estimators, are stated. Section 4 analyzes some specific sieve priors. Section 5 contains the proofs of the main results. Section 6 concludes with some further discussions. Appendix contains some omitted technical results.

**Notation:** For any given (random or non-random, scalar or vector) sequence  $\{z_i\}_{i=1}^n$ , we use the notation  $\mathbb{E}_n[z_i] = n^{-1} \sum_{i=1}^n z_i$ , which should be distinguished from the population expectation  $\mathbb{E}[\cdot]$ . For any vector  $z$ , let  $z^{\otimes 2} = zz^T$  where  $z^T$  is the transpose of  $z$ . For any two sequences of positive constants  $r_n$  and  $s_n$ , we write  $r_n \lesssim s_n$  if the ratio  $r_n/s_n$  is bounded, and  $r_n \sim s_n$  if  $r_n \lesssim s_n$  and  $s_n \lesssim r_n$ . Let  $L_2[0, 1]$  denote the usual  $L_2$  space with respect to the Lebesgue measure for functions defined on  $[0, 1]$ . Let  $\|\cdot\|$  denote the  $L_2$ -norm, i.e.,  $\|f\|^2 = \int_0^1 f^2(x)dx$ . The inner product in  $L_2[0, 1]$  is denoted by  $\langle \cdot, \cdot \rangle$ , i.e.,  $\langle f, g \rangle = \int_0^1 f(x)g(x)dx$ . Let  $C[0, 1]$  denote the metric space of all continuous functions on  $[0, 1]$ , equipped with the uniform metric. The Euclidean norm is denoted by  $\|\cdot\|_{\ell^2}$ . For any matrix  $A$ , let  $s_{\min}(A)$  and  $s_{\max}(A)$  denote the minimum and maximum singular values of  $A$ , respectively. Let  $\|A\|_{\text{op}}$  denote the operator norm of a matrix  $A$  (i.e.,  $\|A\|_{\text{op}} = s_{\max}(A)$ ). Denote by  $dN(\mu, \Sigma)(x)$  the density of the multivariate normal distribution with mean vector  $\mu$  and covariance matrix  $\Sigma$ .

**2. Quasi-Bayesian analysis: informal discussion.** In this section, we outline a quasi-Bayesian analysis of the NPIV model (1). The discussion here is informal. The formal discussion is given in Section 3.

Let  $\mathcal{G}$  be a parameter space (say, some smoothness class of functions, such as a Hölder-Zygmund or Sobolev space), for which we assume  $g_0 \in \mathcal{G}$ . We assume that  $\mathcal{G}$  is at least contained in  $C[0, 1]$ :  $\mathcal{G} \subset C[0, 1]$ . Define the conditional moment function as  $m(W, g) = \mathbb{E}[Y - g(X) | W]$ ,  $g \in \mathcal{G}$ . Then  $g_0$  satisfies the conditional moment restriction

$$(3) \quad m(W, g_0) = 0, a.s.$$

Equivalently, we have  $\mathbb{E}[m^2(W, g_0)] = 0$ .

In this paper, for the purpose of robustness, any specific distribution of  $(Y, X, W)$  is not assumed, which we believe is more practical in statistical and econometric applications. So a Bayesian analysis in the standard sense is not applicable here since a proper likelihood for  $g$  ( $g$  is a generic version of  $g_0$ ) is not available. Instead, we use a quasi-likelihood induced from the conditional moment restriction (3).

Let  $(Y_1, X_1, W_1), \dots, (Y_n, X_n, W_n)$  be i.i.d. observations of  $(Y, X, W)$ . Let  $W^n = \{W_1, \dots, W_n\}$  and  $\mathcal{D}_n = \{(Y_1, X_1, W_1), \dots, (Y_n, X_n, W_n)\}$ . By (3), a plausible candidate of the quasi-likelihood would be

$$p_g(W^n) = \exp\{-(n/2)\mathbb{E}_n[m^2(W_i, g)]\},$$

since  $p_g(W^n)$  is maximized at the true structural function  $g_0$ . However, this  $p_g(W^n)$  is infeasible since  $m(\cdot, g)$  is unknown. Instead of using  $p_g(W^n)$ , we replace  $m(\cdot, g)$  by a suitable estimate  $\hat{m}(\cdot, g)$  and use the quasi-likelihood of the form

$$p_g(\mathcal{D}_n) = \exp\{-(n/2)\mathbb{E}_n[\hat{m}^2(W_i, g)]\}.$$

Below we use a wavelet series estimator of  $m(\cdot, g)$ .

The quasi-Bayesian analysis considered here uses this quasi-likelihood as if it were a proper likelihood and puts priors on  $g \in \mathcal{G}$ . In this paper, as in [46], we shall use sieve priors (more precisely, priors constructed on slowly growing sieves; [46] indeed considered another class of priors, see their supplementary material). The basic idea is to construct a sequence of finite dimensional sieves (say,  $\mathcal{G}_n$ ) that well approximates the parameter space  $\mathcal{G}$  (i.e., each function in  $\mathcal{G}$  is well approximated by some function in  $\mathcal{G}_n$  as  $n$  becomes large), and put priors concentrating on these sieves. Each sieve space is a subset of a linear space spanned by some basis functions. Hence the problem reduces to putting priors on the coefficients on those basis functions. Such priors are typically called “(finite dimensional) sieve priors” (or “series priors”) and have been widely used in the nonparametric Bayesian and quasi-Bayesian analysis [see e.g. 24, 54, 25].

Let  $\Pi_n$  be a so-constructed prior on  $g \in \mathcal{G}$ . Then, formally, the posterior-like distribution of  $g$  given  $\mathcal{D}_n$  may be defined by

$$(4) \quad \Pi_n(dg | \mathcal{D}_n) = \frac{p_g(\mathcal{D}_n)\Pi_n(dg)}{\int p_g(\mathcal{D}_n)\Pi_n(dg)},$$

which we call “quasi-posterior distribution”. The quasi-posterior distribution is not a proper posterior distribution in the strict Bayesian sense since  $p_g(\mathcal{D}_n)$  is not a proper likelihood. Nevertheless,  $\Pi_n(dg | \mathcal{D}_n)$  is a proper distribution, i.e.,  $\int \Pi_n(dg | \mathcal{D}_n) = 1$ . Similarly to proper posterior distributions, contraction of the quasi-posterior distribution around  $g_0$  intuitively means that it contains more and more accurate information about the true structural function  $g_0$  as the sample size increases. Hence, as in proper posterior distributions, it is of fundamental importance to study rates of contraction of quasi-posterior distributions. Here we say that the quasi-posterior  $\Pi_n(dg | \mathcal{D}_n)$  contracts around  $g_0$  at rate  $\varepsilon_n \rightarrow 0$  if  $\Pi_n(g : \|g - g_0\| > \varepsilon_n | \mathcal{D}_n) \xrightarrow{P} 0$ .

This quasi-posterior corresponds to what [63] called ‘‘Gibbs algorithm’’ and what [37] called ‘‘Gibbs posterior’’. The framework of the quasi-posterior (Gibbs posterior) allows us a flexibility since a stringent distributional assumption, such as normality, on the data generating process is not required. Such a framework widens a Bayesian approach to broad fields of statistical problems.<sup>4</sup> Moreover the following proposition gives an interesting interpretation of the quasi-posterior.

PROPOSITION 1. *Let  $\eta > 0$  be a fixed constant. Let  $\Pi$  be a prior distribution for  $g$  defined on, say, the Borel  $\sigma$ -field of  $C[0, 1]$ . Suppose that the data  $\mathcal{D}_n$  are fixed and the maps  $g \mapsto \hat{m}_i(W_i, g)$  are measurable with respect to the Borel  $\sigma$ -field of  $C[0, 1]$ . Then, the distribution*

$$\hat{\Pi}_\eta(dg) = \frac{\exp(-\eta \sum_{i=1}^n \hat{m}_i^2(W_i, g)) \Pi(dg)}{\int \exp(-\eta \sum_{i=1}^n \hat{m}_i^2(W_i, g)) \Pi(dg)},$$

*minimizes the empirical information complexity defined by*

$$(5) \quad \check{\Pi} \mapsto \int \sum_{i=1}^n \hat{m}_i^2(W_i, g) \check{\Pi}(dg) + \eta^{-1} D_{KL}(\check{\Pi} \parallel \Pi)$$

*over all distributions  $\check{\Pi}$  absolutely continuous with respect to  $\Pi$ . Here*

$$D_{KL}(\check{\Pi} \parallel \Pi) = \int \check{\pi} \log \check{\pi} \Pi(dg), \text{ with } d\check{\Pi}/d\Pi = \check{\pi},$$

*is the Kullback-Leibler divergence from  $\check{\Pi}$  to  $\Pi$ .*

PROOF. Immediate from Zhang [62, Proposition 5.1]. □

The proposition shows that, given the data  $\mathcal{D}_n$  and a prior  $\Pi = \Pi_n$  on  $g$ , the quasi-posterior  $\Pi_n(dg \mid \mathcal{D}_n)$  defined in (4) is obtained as a minimizer of the empirical information complexity defined by (5) with  $\eta = 1/2$ . This gives a rationale to use  $\Pi_n(dg \mid \mathcal{D}_n)$  as a quasi-posterior since, among all possible ‘‘quasi-posteriors’’, this  $\Pi_n(dg \mid \mathcal{D}_n)$  optimally balances the average of the natural loss function  $g \mapsto \sum_{i=1}^n \hat{m}_i^2(W_i, g)$  and its complexity (or deviation) relative to the initial prior distribution measured by the Kullback-Leibler divergence. The scaling constant (‘‘temperature’’)  $\eta$  is typically treated as

<sup>4</sup>[37, p.2211] remarked: ‘‘This framework of the Gibbs posterior has been overlooked by most statisticians for a long time [...] a foundation for understanding the statistical behavior of the Gibbs posterior, which we believe will open a productive new line of research.’’

a fixed constant [see, e.g., 63, 37]. An alternative way is to choose  $\eta$  in a data-dependent manner, by e.g. cross validation as mentioned in [63]. It is not difficult to see that the theory below can be extended to the case where  $\eta$  is even random, as long as  $\eta$  converges in probability to a fixed positive constant. However, for the sake of simplicity, we take  $\eta = 1/2$  as a benchmark choice (note that as long as  $\eta$  is a fixed positive constant, the analysis can be reduced to the case with  $\eta = 1/2$  by renormalization).

The quasi-posterior distribution provides point estimators of  $g_0$ . A most natural estimator would be the estimator defined by the posterior expectation (the expectation of the quasi-posterior distribution), i.e.,

$$(6) \quad \hat{g}_{QB} = \begin{cases} \int g \Pi_n(dg | \mathcal{D}_n), & \text{if the right integral exists,} \\ 0, & \text{otherwise,} \end{cases}$$

where the integral  $\int g \Pi_n(dg | \mathcal{D}_n)$  is understood as pointwise.

REMARK 1. Quasi-Bayesian approaches (not necessarily in the present form) are widely used and there are several other attempts of making probabilistic interpretation of such approaches. See for example [39] where the “limited information likelihood” is derived as the “best” (in a suitable sense) approximation to the true likelihood function under a set of moment restrictions and the Bayesian analysis with the limited information likelihood is argued ([46] adapted this approach to conditional moment restriction models), and [53] where a version of the empirical likelihood is interpreted in a Bayesian framework.

**3. Main results.** In this section, we study the asymptotic properties of the quasi-posterior distribution and the quasi-Bayes estimator. In doing so, we have to specify certain regularity properties, such as the smoothness of  $g_0$  and the degree of ill-posedness of the problem. How to characterize the “smoothness” of  $g_0$  is important since it is related to how to put priors. For this purpose, we find wavelet theory useful, and use sieve spaces constructed by using wavelet bases.

3.1. *Posterior construction.* To construct quasi-posterior distributions, we have to estimate  $m(\cdot, g)$  and construct a sequence of sieve spaces for  $\mathcal{G}^s$  on which priors concentrate. For the former purpose, we use a (wavelet) series estimator of  $m(\cdot, g)$ , as in [2] and [11]. For the latter purpose, we construct a sequence of sieve spaces formed by the wavelet basis.

We begin with stating the parameter space for  $g_0$  and the wavelet basis used. We assume that the parameter space  $\mathcal{G}$  is either  $(B_{\infty, \infty}^s, \|\cdot\|_{s, \infty, \infty})$

(Hölder-Zygmund space) or  $(B_{2,2}^s, \|\cdot\|_{s,2,2})$  (Sobolev space), where  $B_{p,q}^s$  is the Besov space of functions on  $[0, 1]$  with parameter  $(s, p, q)$  (the parameter  $s$  generally corresponds to “smoothness”; we add “ $s$ ” on the parameter space,  $\mathcal{G} = \mathcal{G}^s$ , to clarify its dependence on  $s$ ). See Section A.2 for the definition of Besov spaces. We assume that  $s > 1/2$ , under which  $\mathcal{G}^s \subset C[0, 1]$ .

Fix (sufficiently large)  $J_0 \geq 0$ , and let  $\{\varphi_{J_0 k}^{\text{int}}\}_{k=0}^{2^{J_0}-1} \cup \{\psi_{jk}^{\text{int}}, j \geq J_0, k = 0, \dots, 2^j - 1\}$  be an  $S$ -regular Cohen-Daubechies-Vial (CDV) wavelet basis for  $L_2[0, 1]$  ([15]), where  $S$  is a positive integer larger than  $s$ . See Appendix A.1 for CDV wavelet bases. For the notational convenience we write  $\phi_1 = \varphi_{J_0,0}^{\text{int}}, \phi_2 = \varphi_{J_0,1}^{\text{int}}, \dots, \phi_{2^{J_0}} = \varphi_{J_0,2^{J_0}-1}^{\text{int}}$ , and  $\phi_{2^j+1} = \psi_{j,0}^{\text{int}}, \phi_{2^j+2} = \psi_{j,1}^{\text{int}}, \dots, \phi_{2^{j+1}} = \psi_{j,2^j-1}^{\text{int}}$  for  $j \geq J_0$ . Here and in what follows:

Take and fix an  $S$ -regular CDV wavelet basis of  $\{\phi_l, l \geq 1\}$  with  $S > s$ , and we keep this convention. Let  $V_j$  be the linear subspace of  $L_2[0, 1]$  spanned by  $\{\phi_1, \dots, \phi_{2^j}\}$ , and denote by  $P_j$  the projection operator onto  $V_j$ , i.e., for any  $g = \sum_{l=1}^{\infty} b_l \phi_l \in L_2[0, 1]$ ,  $P_j g = \sum_{l=1}^{2^j} b_l \phi_l$ . In what follows, for any  $J \in \mathbb{N}$ , the notation  $b^J$  means that it is a vector of dimension  $2^J$ . For example,  $b^J = (b_1, \dots, b_{2^J})^T$ .

REMARK 2 (Approximation property). For either  $g \in B_{\infty,\infty}^s$  or  $B_{2,2}^s$ , we have  $\|g - P_J g\|^2 \leq C 2^{-2Js}$  for all  $J \geq J_0$ . Here the constant  $C$  depends only on  $s$  and the corresponding Besov norm of  $g$ .

REMARK 3. The use of CDV wavelet bases is not crucial and one may use other reasonable bases such as the Fourier and Hermite polynomial bases. The theory below can be extended to such bases with some modifications. However, CDV wavelet bases are particularly well suited to approximate (not necessarily periodic) smooth functions, which is the reason why we use here CDV wavelet bases. On the other hand, for example, the Fourier basis is only appropriate to approximate periodic functions and it is often not natural to assume that the structural function  $g_0$  is periodic.

We shall now move to the posterior construction. For  $J \geq J_0$ , define the  $2^J$ -dimensional vector of functions  $\phi^J(w)$  by

$$\phi^J(w) = (\phi_1(w), \dots, \phi_{2^J}(w))^T.$$

Let  $J_n \geq J_0$  be a sequence of positive integers such that  $J_n \rightarrow \infty$  and  $2^{J_n} = o(n)$ . Then a wavelet series estimator of  $m(\cdot, g)$  is defined as

$$\hat{m}(w, g) = \phi^{J_n}(w)^T (\mathbb{E}_n[\phi^{J_n}(W_i)^{\otimes 2}])^{-1} \mathbb{E}_n[\phi^{J_n}(W_i)(Y_i - g(X_i))],$$

where we replace the inverse matrix by the generalized inverse if the former does not exist; the probability of such an event converges to zero as  $n \rightarrow \infty$  under the assumptions below. We use this wavelet series estimator throughout the analysis.

For the same  $J_n$ , we shall take  $V_{J_n} = \text{span}\{\phi_1, \dots, \phi_{2^{J_n}}\}$  as a sieve space for  $\mathcal{G}^s$ . We consider priors  $\Pi_n$  that concentrate on  $V_{J_n}$ , i.e.,  $\Pi_n(V_{J_n}) = 1$ . Formally, we think of that priors on  $g$  are defined on the Borel  $\sigma$ -field of  $C[0, 1]$  (hence the quasi-posterior  $\Pi_n(dg | \mathcal{D}_n)$  is understood to be defined on the Borel  $\sigma$ -field of  $C[0, 1]$ , which is possible since the map  $g \mapsto p_g(\mathcal{D}_n)$  is continuous on  $C[0, 1]$ ). Since the map  $b^{J_n} = (b_1, \dots, b_{2^{J_n}})^T \mapsto \sum_{l=1}^{2^{J_n}} b_l \phi_l, \mathbb{R}^{2^{J_n}} \rightarrow C[0, 1]$ , is homeomorphic from  $\mathbb{R}^{2^{J_n}}$  onto  $V_{J_n}$ , putting priors on  $g \in V_{J_n}$  is equivalent to putting priors on  $b^{J_n} \in \mathbb{R}^{2^{J_n}}$  (the latter are of course defined on the Borel  $\sigma$ -field of  $\mathbb{R}^{2^{J_n}}$ ). Practically, priors on  $g \in V_{J_n}$  are induced from priors on  $b^{J_n} \in \mathbb{R}^{2^{J_n}}$ . For the later purpose, it is useful to determine the correspondence between priors for these two parameterizations. Unless otherwise stated, we follow the convention of the notation such that:

$$\tilde{\Pi}_n: \text{a prior on } b^{J_n} \in \mathbb{R}^{2^{J_n}} \leftrightarrow \Pi_n: \text{the induced prior on } g \in V_{J_n}.$$

We shall call  $\tilde{\Pi}_n$  a generating prior, and  $\Pi_n$  the induced prior.

Correspondingly, the quasi-posterior for  $b^{J_n}$  is defined. With a slight abuse of notation, for  $g = \sum_{l=1}^{2^{J_n}} b_l \phi_l$ , we write  $\hat{m}(w, b^{J_n}) = \hat{m}(w, g)$ , and take  $p_{b^{J_n}}(\mathcal{D}_n) = \exp\{-(n/2)\mathbb{E}_n[\hat{m}^2(W_i, b^{J_n})]\}$  as a quasi-likelihood for  $b^{J_n}$ . Note that in this particular setting, the log quasi-likelihood is quadratic in  $b^{J_n}$ . Let  $\tilde{\Pi}_n(db^{J_n} | \mathcal{D}_n)$  denote the resulting quasi-posterior distribution for  $b^{J_n}$ :

$$(7) \quad \tilde{\Pi}_n(db^{J_n} | \mathcal{D}_n) = \frac{p_{b^{J_n}}(\mathcal{D}_n)\tilde{\Pi}_n(db^{J_n})}{\int p_{b^{J_n}}(\mathcal{D}_n)\tilde{\Pi}_n(db^{J_n})}.$$

For the quasi-Bayes estimator  $\hat{g}_{QB}$  defined by (6), since for every  $x \in [0, 1]$ , the map  $g \mapsto g(x)$  is continuous on  $C[0, 1]$ , and conditional on  $\mathcal{D}_n$  the quasi-posterior  $\Pi_n(dg | \mathcal{D}_n)$  is a Borel probability measure on  $C[0, 1]$ , the integral  $\int g(x)\Pi_n(dg | \mathcal{D}_n)$  exists as soon as  $\int |g(x)|\Pi_n(dg | \mathcal{D}_n) < \infty$ . Furthermore,  $\hat{g}_{QB}$  can be computed by using the relation

$$\int g(x)\Pi_n(dg | \mathcal{D}_n) = \phi^{J_n}(x)^T \left[ \int b^{J_n}\tilde{\Pi}_n(db^{J_n} | \mathcal{D}_n) \right],$$

as soon as the integral on the right side exists. Hence, practically, it is sufficient to compute the expectation of  $\tilde{\Pi}_n(db^{J_n} | \mathcal{D}_n)$ .

REMARK 4. The use of the same wavelet basis to estimate  $m(\cdot, g)$  and to construct a sequence of sieve spaces for  $\mathcal{G}^s$  is not essential and can be

relaxed. Suppose that we have another CDV wavelet basis  $\{\tilde{\phi}_l\}$  for  $L_2[0, 1]$  and use this basis to estimate  $m(\cdot, g)$ . Then, all the results below apply by simply replacing  $\phi_l(W_i)$  by  $\tilde{\phi}_l(W_i)$ . To keep the notation simple, we use the same wavelet basis.

However, the use of the same resolution level  $J_n$  is essential (at least at the proof level) in establishing the asymptotic properties of the quasi-posterior distribution. It may be a technical artifact, but we do not extend the theory in this direction since there is no clear theoretical benefit to do so (note that in the purely frequentist estimation case, [11] allowed for using different cut-off levels for approximating  $m(\cdot, g)$  and  $g(\cdot)$ ).

*3.2. Basic assumptions.* We state some basic assumptions. We do not state here assumptions on priors, which will be stated in the theorems below. In what follows, let  $C_1 > 1$  be a sufficiently large constant.

ASSUMPTION 1. (i)  $(X, W)$  has a joint density  $f_{X,W}(x, w)$  on  $[0, 1]^2$  satisfying that  $f_{X,W}(x, w) \leq C_1, \forall x, w \in [0, 1]$ . (ii)  $\sup_{w \in [0, 1]} \mathbb{E}[U^2 | W = w] \leq C_1$  where  $U = Y - g_0(X)$ . (iii)  $s_{\min}(\mathbb{E}[\phi^J(W)^{\otimes 2}]) \geq C_1^{-1}, \forall J \geq J_0$ .

Assumption 1 is a usual restriction in the literature, up to minor differences [see 28, 34]. Denote by  $f_X(x)$  and  $f_W(w)$  the marginal densities of  $X$  and  $W$ , respectively i.e.,  $f_X(x) = \int f_{X,W}(x, w)dw$  and  $f_W(w) = \int f_{X,W}(x, w)dx$ . Then Assumption 1 (i) implies that  $f_X(x) \leq C_1, \forall x \in [0, 1]$  and  $f_W(w) \leq C_1, \forall w \in [0, 1]$ . A primitive regularity condition that guarantees Assumption 1 (iii) is that  $f_W(w) \geq C_1^{-1}$  for all  $w \in [0, 1]$ . To see this, for  $\alpha^J \in \mathbb{R}^{2^J}$  with  $\|\alpha^J\|_{\ell^2} = 1$ , we have

$$\begin{aligned} (\alpha^J)^T \mathbb{E}[\phi^J(W)^{\otimes 2}] \alpha^J &= \int_0^1 (\phi^J(w)^T \alpha^J)^2 f_W(w) dw \geq C_1^{-1} \int_0^1 (\phi^J(w)^T \alpha^J)^2 dw \\ (8) \quad &= C_1^{-1} (\alpha^J)^T \left[ \int_0^1 \phi^J(w) \phi^J(w)^T dw \right] \alpha^J = C_1^{-1} \|\alpha^J\|_{\ell^2}^2 = C_1^{-1}, \end{aligned}$$

where we have used the fact that  $\{\phi_l\}$  is orthonormal in  $L_2[0, 1]$

For identification of  $g_0$ , we assume:

ASSUMPTION 2. The linear operator  $K : L_2[0, 1] \rightarrow L_2[0, 1]$  is injective.

For smoothness of  $g_0$ , as mentioned before, we assume:

ASSUMPTION 3.  $\exists s > 1/2, g_0 \in \mathcal{G}^s$ , where  $\mathcal{G}^s$  is either  $B_{\infty, \infty}^s$  or  $B_{2, 2}^s$ .

The identification condition (Assumption 2) is equivalent to the “completeness” of the conditional distribution of  $X$  conditional on  $W$  [50]. We refer the reader to [55], [30] and [36] for discussion on the completeness condition. We should note that restricting the domain of  $K$  to a “small” set, such as a Sobolev ball, would substantially relax Assumption 2, which however requires a different analysis. For the sake of simplicity, we assume the injectivity of  $K$  on the full domain.

As discussed in Introduction, solving (2) is an ill-posed inverse problem. Thus, the statistical difficulty of estimating  $g_0$  depends on the difficulty of continuously inverting  $K$ , which is usually referred to as “ill-posedness” of the inverse problem (2). Typically, the ill-posedness is characterized by the decay rate of  $\kappa_l \rightarrow 0$  ( $\kappa_l$  is the  $l$ -th largest singular value of  $K$ ), which is plausible if  $K$  were known and the singular value decomposition of  $K$  were used [see 10]. However, here,  $K$  is unknown and the known wavelet basis  $\{\phi_l\}$  is used instead of the singular value system. Thus, it is suitable to quantify the ill-posedness using the wavelet basis  $\{\phi_l\}$ . To this end, define

$$\tau_J = s_{\min}(\mathbb{E}[\phi^J(W)\phi^J(X)^T]) = s_{\min}(\langle\langle\phi_l, K\phi_m\rangle\rangle_{1 \leq l, m \leq 2^J}), \quad J \geq J_0.$$

This quantity corresponds to (the reciprocal of) what is called “sieve measure of ill-posedness” in the literature [6, 34]. We at least have to assume that  $\tau_J > 0$  for all  $J \geq J_0$ . Note however that

$$\begin{aligned} \tau_J &= s_{\min}(\langle\langle\phi_l, K\phi_m\rangle\rangle_{1 \leq l, m \leq 2^J}) \\ &= \min_{g \in V_J, \|g\|=1} \|\langle\langle\phi_l, Kg\rangle\rangle_{1 \leq l \leq 2^J}\|_{\ell^2} \\ &\leq \min_{g \in V_J, \|g\|=1} \|Kg\| \quad (\text{Bessel's inequality}) \\ &\leq \kappa_{2^J}, \quad (\text{Courant-Fischer-Weyl's minimax principle}) \end{aligned}$$

by which, necessarily,  $\tau_J \rightarrow 0$  as  $J \rightarrow \infty$ . For this quantity, we assume:

ASSUMPTION 4. (i) (*mildly ill-posed case*)  $\exists r > 0$ ,  $\tau_J \geq C_1^{-1}2^{-Jr}$ ,  $\forall J \geq J_0$  or (*severely ill-posed case*)  $\exists c > 0$ ,  $\tau_J \geq C_1^{-1}\exp(-c2^J)$ ,  $\forall J \geq J_0$ ; (ii)

$$\begin{aligned} \|\mathbb{E}[\phi^J(W)(g_0 - P_J g_0)(X)]\|_{\ell^2} & (= \|\langle\langle\phi_l, K(g_0 - P_J g_0)\rangle\rangle_{l=1}^{2^J}\|_{\ell^2}) \\ & \leq C_1 \tau_J \|g_0 - P_J g_0\|, \quad \forall J \geq J_0. \end{aligned}$$

Assumption 4 (i) lower bounds  $\tau_J$  as  $J \rightarrow \infty$ , thereby quantifies the ill-posedness. We cover both the “mildly ill-posed” and “severely ill-posed” cases [this definition of mild ill-posedness and severe ill-posedness is due to

33, 34]. The severely ill-posed case happens, e.g., when the joint density  $f_{X,W}(x, w)$  is analytic [see 45, Theorem 15.20].

Assumption 4 (ii) is a “stability” condition about the bias  $g_0 - P_J g_0$ , which states that  $K(g_0 - P_J g_0)$  is sufficiently “small” relative to  $g_0 - P_J g_0$ . Note that in the (ideal) case in which, e.g.,  $K$  is self-adjoint and  $\{\phi_l\}$  is the eigen-basis of  $K$ ,  $\langle \phi_l, K(g_0 - P_J g_0) \rangle = 0$  for all  $l = 1, \dots, 2^J$ , in which case Assumption 4 (ii) is trivially satisfied. Assumption 4 (ii) allows more general situations in which  $K$  may not be self-adjoint and  $\{\phi_l\}$  may not be the eigen-basis of  $K$  by allowing for a certain “slack”. This assumption, although looks technical, is common in the study of rates of convergence in estimation of the structural function  $g_0$ . Indeed, essentially similar conditions have appeared in the past literature such as [6, 12, 34]. For example, Blundell et al. [6, Assumption 6] essentially states (in our notation) that  $\|K(g_0 - P_J g_0)\| \leq C_1 \tau_J \|g_0 - P_J g_0\|$ , which implies our Assumption 4 (ii) since  $\|(\langle \phi_l, K(g_0 - P_J g_0) \rangle)_{l=1}^{2^J}\|_{\ell^2} \leq \|K(g_0 - P_J g_0)\|$  (Bessel’s inequality).

REMARK 5. For given values of  $C_1 > 1$ ,  $M > 0$ ,  $r > 0$ ,  $c > 0$  and  $s > 1/2$ , let  $\mathcal{F} = \mathcal{F}(C_1, M, r, c, s)$  denote the set of all distributions of  $(Y, X, W)$  satisfying Assumptions 1-4 with  $\|g_0\|_{s, \infty, \infty} \leq M$  in case of  $\mathcal{G}^s = B_{\infty, \infty}^s$  and  $\|g_0\|_{s, 2, 2} \leq M$  in case of  $\mathcal{G}^s = B_{2, 2}^s$ . By [28, 12], it is shown that the minimax rate of convergence (in  $\|\cdot\|$ ) of estimation of  $g_0$  over this distribution class  $\mathcal{F}$  is  $n^{-s/(2r+2s+1)}$  in the mildly ill-posed case (where  $\tau_J \geq C_1^{-1} 2^{-Jr}$ ) and  $(\log n)^{-s}$  in the severely ill-posed case (where  $\tau_J \geq C_1^{-1} \exp(-c2^J)$ ) as the sample size  $n \rightarrow \infty$  (the assumption on the conditional second moment of  $U$  given  $W$  is not binding; i.e., replacing Assumption 1(ii) by a stronger one, such as  $\sup_{w \in [0, 1]} \mathbb{E}[|U|^{2+\epsilon} | W = w] \leq C_1$  for some  $\epsilon > 0$  determined outside the class of distributions, does not alter these minimax rates).

By Theorem 2.5 of [24], it is readily seen that these rates are the fastest possible rates of contraction of (general) quasi-posterior distributions in this setting. More formally, we can state the following assertion:

*Let  $\Pi_n(dg | \mathcal{D}_n)$  be the quasi-posterior distribution defined on, say, the Borel  $\sigma$ -field of  $C[0, 1]$ , constructed from putting a suitable prior on  $g$  to the quasi-likelihood  $p_g(\mathcal{D}_n)$  (the prior here needs not be a sieve prior). Suppose now that for some  $\varepsilon_n \rightarrow 0$ ,  $\sup_{F \in \mathcal{F}} \mathbb{E}_F[\Pi_n(g : \|g - g_0\| > \varepsilon_n | \mathcal{D}_n)] \rightarrow 0$ . Then there exists a point estimator that converges (in probability) at least as fast as  $\varepsilon_n$  uniformly in  $F \in \mathcal{F}$ .*

The proof is just a small modification of that of Theorem 2.5 in [24] and hence omitted. Importantly, the quasi-posterior cannot contract at a rate faster than the optimal rate of convergence for point estimators [24, page 507,

lines 19-20]. Hence, in the minimax sense, the fastest possible rate of contraction of the quasi-posterior distribution  $\Pi_n(dg \mid \mathcal{D}_n)$  is  $n^{-s/(2r+2s+1)}$  in the mildly ill-posed case and  $(\log n)^{-s}$  in the severely ill-posed case (Proposition 2 in Section 4 ahead shows that these rates are indeed attainable for suitable sieve priors).

**3.3. Main results: general theorems.** This section presents general theorems on contraction rates and asymptotic normality for quasi-posterior distributions as well as convergence rates for quasi-Bayes estimators. In what follows, let  $(Y_1, X_1, W_1), \dots, (Y_n, X_n, W_n)$  be i.i.d. observations of  $(Y, X, W)$ . Denote by  $b_0^J = (b_{01}, \dots, b_{0,2^J})^T$  the vector of the first  $2^J$  generalized Fourier coefficients of  $g_0$ , i.e.,  $b_{0l} = \int \phi_l g_0$ . Let  $\|\cdot\|_{\text{TV}}$  denote the total variation norm between two distributions.

**THEOREM 1.** *Suppose that Assumptions 1-4 are satisfied. Take  $J_n$  in such a way that  $J_n \rightarrow \infty$  and  $J_n 2^{J_n}/n = o(\tau_{J_n}^2)$ . Let  $\epsilon_n$  be a sequence of positive constants such that  $\epsilon_n \rightarrow 0$  and  $n\epsilon_n^2 \gtrsim 2^{J_n}$ . Suppose that generating priors  $\tilde{\Pi}_n$  has densities  $\tilde{\pi}_n$  on  $\mathbb{R}^{2^{J_n}}$  and satisfy the following conditions:*

- P1) (Small ball condition) There exists a constant  $C > 0$  such that for all  $n$  sufficiently large,  $\tilde{\Pi}_n(b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\ell^2} \leq \epsilon_n) \geq e^{-Cn\epsilon_n^2}$ .*  
*P2) (Prior flatness condition) Let  $\gamma_n = 2^{-J_n s} + \tau_{J_n}^{-1}\epsilon_n$ . There exists a sequence of constants  $L_n \rightarrow \infty$  sufficiently slowly such that for all  $n$  sufficiently large,  $\tilde{\pi}_n(b^{J_n})$  is positive for all  $\|b^{J_n} - b_0^{J_n}\|_{\ell^2} \leq L_n \gamma_n$ , and*

$$\sup_{\|b^{J_n}\|_{\ell^2} \leq L_n \gamma_n, \|\tilde{b}^{J_n}\|_{\ell^2} \leq L_n \gamma_n} \left| \frac{\tilde{\pi}_n(b_0^{J_n} + b^{J_n})}{\tilde{\pi}_n(b_0^{J_n} + \tilde{b}^{J_n})} - 1 \right| \rightarrow 0.$$

Then for every sequence  $M_n \rightarrow \infty$ , we have

$$(9) \quad \tilde{\Pi}_n \left\{ b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\ell^2} > M_n (2^{-J_n s} + \tau_{J_n}^{-1} \sqrt{2^{J_n}/n}) \mid \mathcal{D}_n \right\} \xrightarrow{P} 0.$$

Furthermore, assume that  $J_n 2^{3J_n}/n = o(\tau_{J_n}^2)$ . Then we have

$$\|\tilde{\Pi}_n(\cdot \mid \mathcal{D}_n) - N(\hat{b}^{J_n}, n^{-1} \Phi_{WX}^{-1} \Phi_{WW} \Phi_{XW}^{-1})(\cdot)\|_{\text{TV}} \xrightarrow{P} 0,$$

where  $\Phi_{WX} := \mathbb{E}[\phi^{J_n}(W)\phi^{J_n}(X)^T]$ ,  $\Phi_{XW} := \Phi_{WX}^T$ ,  $\Phi_{WW} := \mathbb{E}[\phi^{J_n}(W)^{\otimes 2}]$ , and where  $\hat{b}^{J_n}$  is a “maximum quasi-likelihood estimator” of  $b_0^{J_n}$ , i.e.,

$$(10) \quad \hat{b}^{J_n} \in \arg \max_{b^{J_n} \in \mathbb{R}^{2^{J_n}}} p_{b^{J_n}}(\mathcal{D}_n).$$

PROOF. See Section 5.1.  $\square$

REMARK 6. The condition  $J_n 2^{J_n}/n = o(\tau_{J_n}^2)$  appears essentially because the operator  $K$  is unknown. In our setup, this results in estimating the matrix  $\mathbb{E}[\phi^{J_n}(W)\phi^{J_n}(X)^T]$  by its empirical counterpart  $\mathbb{E}_n[\phi^{J_n}(W_i)\phi^{J_n}(X_i)^T]$ . In the proof, we have to suitably lower bound the minimum singular value of  $\mathbb{E}_n[\phi^{J_n}(W_i)\phi^{J_n}(X_i)^T]$ , denoted by  $\hat{\tau}_{J_n}$ , which is an empirical counterpart of the sieve measure of ill-posedness  $\tau_{J_n}$ . By Lemma 1, we have  $\hat{\tau}_{J_n} = \tau_{J_n} - O_P(\sqrt{J_n 2^{J_n}/n})$ , so that to make the estimation effect in  $\hat{\tau}_{J_n}$  negligible, we need  $J_n 2^{J_n}/n = o(\tau_{J_n}^2)$ .

REMARK 7. Theorem 1 is abstract in the sense that it only gives conditions P1) and P2) on priors for which (9) and (10) hold. For specific priors, we have to check these conditions with possible  $J_n$ , which will be done in Section 4.

Since for  $g = \sum_{l=1}^{2^{J_n}} b_l \phi_l$ ,  $\|g - g_0\|^2 = \|g - P_{J_n} g_0\|^2 + \|g_0 - P_{J_n} g_0\|^2 \lesssim \|b^{J_n} - b_0^{J_n}\|_{\ell^2}^2 + 2^{-2J_n s}$ , part (9) of Theorem 1 leads to that for every sequence  $M_n \rightarrow \infty$ , we have

$$\Pi_n \left\{ g : \|g - g_0\| > M_n (2^{-J_n s} + \tau_{J_n}^{-1} \sqrt{2^{J_n}/n}) \mid \mathcal{D}_n \right\} \xrightarrow{P} 0,$$

which means that the rate of contraction of the quasi-posterior distribution  $\Pi_n(dg \mid \mathcal{D}_n)$  is  $\max\{2^{-J_n s}, \tau_{J_n}^{-1} \sqrt{2^{J_n}/n}\}$ .<sup>5</sup> In many examples, for given  $J_n \rightarrow \infty$  with  $J_n 2^{J_n}/n = o(\tau_{J_n}^2)$ , condition P1) is satisfied with  $\epsilon_n \sim \sqrt{2^{J_n}(\log n)/n}$ . Taking  $J_n$  in such a way that (with some constant  $c' < 1/(2c)$  in the severely ill-posed case)

$$(11) \quad \begin{cases} 2^{J_n} \sim n^{1/(2r+2s+1)}, & \text{in the mildly ill-posed case,} \\ \lim_{n \rightarrow \infty} (2^{J_n}/(c' \log n)) = 1, & \text{in the severely ill-posed case,} \end{cases}$$

under which the optimal contraction rate is attained,  $\gamma_n$  in condition P2) is

$$(12) \quad \gamma_n \sim \begin{cases} n^{-s/(2r+2s+1)} (\log n)^{1/2}, & \text{in the mildly ill-posed case,} \\ (\log n)^{-s}, & \text{in the severely ill-posed case,} \end{cases}$$

So condition P2) states that, to attain the optimal contraction rate (and the Bernstein-von Mises type result), the prior density  $\tilde{\pi}_n$  should be sufficiently

<sup>5</sup>We have ignored the appearance of  $M_n \rightarrow \infty$ , which can be arbitrarily slow. A version in which  $M_n$  is replaced by a large fixed constant  $M > 0$  is presented in Theorem 2.

“flat” in a ball with center  $b_0^{J_n}$  and radius of order (12). Some specific priors leading to the optimal contraction rate will be given in Section 4.

As noted before, in many examples, for given  $J_n \rightarrow \infty$  with  $J_n 2^{J_n}/n = o(\tau_{J_n}^2)$ , condition P1) is satisfied with  $\epsilon_n \sim \sqrt{2^{J_n}(\log n)/n}$ . Inspection of the proof shows that, without condition P2), this already leads to contraction rate  $\max\{2^{-J_n s}, \tau_{J_n}^{-1} \sqrt{2^{J_n}(\log n)/n}\}$ , which, in the mildly ill-posed case, reduces to  $(n/\log n)^{-s/(2r+2s+1)}$  by taking  $2^{J_n} \sim (n/\log n)^{1/(2r+2s+1)}$ . However, this rate is not fully satisfactory because of the appearance of the log term. Condition P2) is used to get rid of the log term.

The small ball condition P1) is standard in nonparametric Bayesian statistics and analogous to condition (2.4) in [24]. It is however stated in Ghosal et al. [24, p.505-506] that their Theorem 2.1 is not sharp enough when priors constructed on a sequence of finite dimensional sieves are used, and more sophisticated condition (2.9) is devised in their Theorem 2.4 (see also the proof of their Theorem 4.5). However, a version of their condition (2.9) is not clear to work in our problem, because the effect of the random matrix  $\mathbb{E}_n[\phi^{J_n}(W_i)\phi^{J_n}(X_i)^T]$  has to be suitably controlled. Instead, we devise condition P2) to obtain sharper contraction rates.

Under a further integrability condition about  $U = Y - g_0(X)$ ,  $M_n \rightarrow \infty$  in (9) can be replaced by a large fixed constant  $M$ .

**THEOREM 2.** *Suppose that all the conditions that guarantee (9) in Theorem 1 are satisfied. Furthermore, assume that  $\sup_{w \in [0,1]} \mathbb{E}[U^2 1(|U| > \lambda) | W = w] \rightarrow 0$  as  $\lambda \rightarrow \infty$  where  $U = Y - g_0(X)$ . Then there exists a constant  $M > 0$  such that*

$$(13) \quad \tilde{\Pi}_n \left\{ b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\ell^2} > M(2^{-J_n s} + \tau_{J_n}^{-1} \sqrt{2^{J_n}/n}) \mid \mathcal{D}_n \right\} \xrightarrow{P} 0.$$

**PROOF.** See Section 5.2. □

The proof consists in establishing a concentration property of the random variable  $\|\mathbb{E}_n[\phi^{J_n}(W_i)U_i]\|_{\ell^2}$ , which uses a truncation argument and Talagrand’s [58] concentration inequality. A sufficient condition that guarantees that  $\sup_{w \in [0,1]} \mathbb{E}[U^2 1(|U| > \lambda) | W = w] \rightarrow 0$  as  $\lambda \rightarrow \infty$  is that  $\exists \epsilon > 0$ ,  $\sup_{w \in [0,1]} \mathbb{E}[|U|^{2+\epsilon} | W = w] < \infty$ . The additional condition in Theorem 2 is a uniform integrability condition and stronger than Assumption 1 (ii). To see this, note that  $U$  is distributed as  $F_{U|W}^{-1}(\mathcal{U} | W)$  where  $F_{U|W}^{-1}(u | w)$  is the conditional quantile function of  $U$  given  $W = w$ , and  $\mathcal{U}$  is a uniform random variable on  $(0, 1)$  independent of  $W$ . Think of  $U_w(u) = F_{U|W}^{-1}(u | w)$ ,  $w \in [0, 1]$  as a stochastic process defined on  $((0, 1), \mu)$  with  $\mu$  Lebesgue measure

on  $(0, 1)$ . Then the condition  $\sup_{w \in [0, 1]} \mathbb{E}[U^2 1(|U| > \lambda) \mid W = w] \rightarrow 0$  as  $\lambda \rightarrow \infty$  states exactly the uniform integrability of  $(U_w)_{w \in [0, 1]}$ .

The second part of Theorem 1 states a Bernstein-von Mises type result for the quasi-posterior distribution  $\tilde{\Pi}_n(db^{J_n} \mid \mathcal{D}_n)$ , which states that the quasi-posterior distribution is approximated by the normal distribution centered at  $\hat{b}^{J_n}$ , which is often referred to as the ‘‘sieve minimum distance estimator’’ and is a benchmark frequentist estimator for these types of models. Note that, neglecting the bias,  $\hat{b}^{J_n}$  is approximated as  $b_0^{J_n} + \Phi_{WX}^{-1} \mathbb{E}_n[\phi^{J_n}(W_i)U_i]$ , but the covariance matrix of the term  $\Phi_{WX}^{-1} \sqrt{n} \mathbb{E}_n[\phi^{J_n}(W_i)U_i]$  is generally different from  $\Phi_{WX}^{-1} \Phi_{WW} \Phi_{XW}^{-1}$  (which is the reason why we added ‘‘type’’). This is a generic nature of quasi-posterior distributions. Even for finite dimensional models, generally, the covariance matrix of the centering variable does not coincide with that of the normal distribution approximating the quasi-posterior distribution [see 13].

Lastly, we consider the convergence rate of the quasi-Bayes estimator  $\hat{g}_{QB}$  of  $g_0$  defined by (6).

**THEOREM 3.** *Suppose that all the conditions of Theorem 2 are satisfied. Let  $\hat{g}_{QB}$  be the quasi-Bayes estimator defined by (6). Then  $\mathbb{P}\{\mathcal{D}_n : \int |g(x)| \Pi_n(dg \mid \mathcal{D}_n) < \infty, \forall x \in [0, 1]\} \rightarrow 1$ , and there exists a constant  $D > 0$  such that for every sequence  $M_n \rightarrow \infty$ ,*

$$(14) \quad \mathbb{P}\left[\|\hat{g}_{QB} - g_0\| \leq D \max\{2^{-J_n s}, \tau_{J_n}^{-1} \sqrt{2^{J_n}/n}, \tau_{J_n}^{-1} \epsilon_n \varrho_n M_n\}\right] \rightarrow 1,$$

where

$$\varrho_n := \sup_{\|b^{J_n}\|_{\ell_2} \leq L_n \gamma_n, \|\tilde{b}^{J_n}\|_{\ell_2} \leq L_n \gamma_n} \left| \frac{\tilde{\pi}_n(b_0^{J_n} + b^{J_n})}{\tilde{\pi}_n(b_0^{J_n} + \tilde{b}^{J_n})} - 1 \right|,$$

and where  $\epsilon_n, \gamma_n$  and  $L_n$  are given in the statement of Theorem 1.

**PROOF.** See Appendix C. □

Theorem 3 is *not* directly deduced from Theorem 1. Indeed,  $\|g - g_0\|$  may be unbounded on the support of  $\Pi_n$  since the support of  $\Pi_n$  may be unbounded in  $\|\cdot\|$ , and hence the argument in Ghosal et al. [24, p.506-p.507] can not apply (in [24], a typical distance to measure the goodness of a point estimator is the Hellinger distance and uniformly bounded). Hence an additional work is needed to prove Theorem 3.

The convergence rate of the quasi-Bayes estimator is determined by the three terms:  $2^{-J_n s}$ ,  $\tau_{J_n}^{-1} \sqrt{2^{J_n}/n}$ , and  $\tau_{J_n}^{-1} \epsilon_n \varrho_n M_n$ . The last term is typically small relative to the other two terms. Indeed, as noted before, in

many examples, for given  $J_n \rightarrow \infty$  with  $J_n 2^{J_n}/n = o(\tau_{J_n}^2)$ ,  $\epsilon_n$  can be taken in such a way that  $\epsilon_n \sim \sqrt{2^{J_n}(\log n)/n}$ . In that case  $\tau_{J_n}^{-1} \epsilon_n \varrho_n M_n \sim \tau_{J_n}^{-1} \varrho_n M_n \sqrt{2^{J_n}(\log n)/n}$ , and as long as  $\varrho_n \rightarrow 0$  sufficiently fast, i.e.,  $\varrho_n = o((\log n)^{-1/2})$ , the convergence rate of the quasi-Bayes estimator  $\hat{g}_{QB}$  reduces to  $\max\{2^{-J_n s}, \tau_{J_n}^{-1} \sqrt{2^{J_n}/n}\}$ .

**4. Prior specification: examples.** In this section, we give some specific sieve priors for which the quasi-posterior distribution (the quasi-Bayes estimator) attains the minimax optimal rate of contraction (convergence, respectively). We consider two types of priors, namely, product and isotropic priors. We will verify that these priors meet conditions P1) and P2) in Theorem 1 with the choice (11). For the notational convenience, define

$$\epsilon_{n,s,r} = \begin{cases} n^{-s/(2s+2r+1)}, & \text{in the mildly ill-posed case,} \\ (\log n)^{-s}, & \text{in the severely ill-posed case.} \end{cases}$$

We may think of the severely ill-posed case as the case with  $r = \infty$ .

**PROPOSITION 2.** *Suppose that Assumptions 1-4 are satisfied. Consider the following two classes of prior distributions on  $\mathbb{R}^{2^{J_n}}$ :*

**(Product prior)** *Let  $q(x)$  be a probability density function on  $\mathbb{R}$  such that for a constant  $A > \sup_{l \geq 1} |b_{0l}|$ : 1)  $q(x)$  is positive on  $[-A, A]$ ; 2)  $\log q(x)$  is Lipschitz continuous on  $[-A, A]$ , i.e., there exists a constant  $L > 0$  possibly depending on  $A$  such that  $|\log q(x) - \log q(y)| \leq L|x - y|$ ,  $\forall x, y \in [-A, A]$ . Take the density of the generating prior by  $\tilde{\pi}_n(b^{J_n}) = \prod_{l=1}^{2^{J_n}} q(b_l)$ .*

**(Isotropic prior)** *Let  $r(x)$  be a probability density function on  $[0, \infty)$  having all moments such that: 1) for a constant  $A > \|g_0\|$ ,  $r(x)$  is positive and continuous on  $[0, A]$ ; 2) for a constant  $c'' > 0$ ,  $\int_0^\infty x^{k-1} r(x) dx \leq e^{c'' k \log k}$  for all  $k$  sufficiently large. Take the density of the generating prior by  $\tilde{\pi}_n(b^{J_n}) \propto r(\|b^{J_n}\|_{\ell^2})$ .*

*Take  $J_n$  as in (11). Then, in either case of product or isotropic priors, for every sequence  $M_n \rightarrow \infty$ , we have  $\Pi_n\{g : \|g - g_0\| > M_n \epsilon_{n,s,r} \mid \mathcal{D}_n\} \xrightarrow{P} 0$ . Furthermore, if  $\sup_{w \in [0,1]} \mathbb{E}[U^2 1(|U| > \lambda) \mid W = w] \rightarrow 0$  as  $\lambda \rightarrow \infty$ , then there exists a constant  $M > 0$  such that  $\Pi_n\{g : \|g - g_0\| > M \epsilon_{n,s,r} \mid \mathcal{D}_n\} \xrightarrow{P} 0$ .*

**PROOF.** See Appendix D. □

Proposition 2 shows that a wide class of priors constructed on slowly growing sieves lead to the minimax optimal contraction rate (see Remark

5). In either case of product or isotropic priors, the constant  $A$  is not necessarily known, which allows  $q(x)$  and  $r(x)$  to have unbounded support. For example, in the former case,  $q(x)$  may be the density of the standard normal distribution, in which case  $A$  can be taken to be arbitrarily large. Likewise, in the latter case,  $r(x)$  may be the density of an exponential distribution:  $r(x) = \lambda e^{-\lambda x}$ ,  $x \geq 0$  for some  $\lambda > 0$ . In the isotropic prior case,  $r(x)$  should have all moments, i.e.,  $\int_0^\infty x^k r(x) dx < \infty$  for all  $k \geq 1$ , which ensures that  $\tilde{\pi}_n(b^{J_n}) \propto r(\|b^{J_n}\|_{\ell^2})$  is a proper distribution on  $\mathbb{R}^{2^{J_n}}$  for every  $n \geq 1$ .

The next proposition shows that two classes of priors in Proposition 2 lead to the minimax optimal convergence rate for the quasi-Bayes estimator.

**PROPOSITION 3.** *Suppose that Assumptions 1-4 are satisfied. Furthermore, assume that  $\sup_{w \in [0,1]} \mathbb{E}[U^2 1(|U| > \lambda) | W = w] \rightarrow 0$  as  $\lambda \rightarrow \infty$ . Consider the two classes of prior distributions on  $\mathbb{R}^{2^{J_n}}$  given in Proposition 2. In the isotropic prior case, assume further that  $r(x)$  is Lipschitz continuous on  $[0, A]$ . Take  $J_n$  as in (11). Then, in either case of product or isotropic priors, there exists a constant  $M > 0$  such that  $\mathbb{P}\{\|\hat{g}_{QB} - g_0\| > M\varepsilon_{n,s,r}\} \rightarrow 0$ .*

**PROOF.** See Appendix D. □

**REMARK 8.** In the above propositions,  $J_n$  plays the role of regularization and should be chosen sufficiently slowly growing, thereby there is no need to place restrictions on weights on  $b_l$  between  $1 \leq l \leq 2^{J_n}$ . The abstract Theorem 1 is derived to cover this case. There is another way to deal with the ill-posedness, i.e., allowing for large dimensional sieves but placing prior distributions that have smaller weights on  $b_l$  for larger  $l$  (“shrinking priors”), which corresponds to the “sieve method using large dimensional sieves with heavy penalties” in the classification of [11].<sup>6</sup> The supplementary material of [46] is concerned with this approach, but they did not establish sharp contraction rates. The extension to this approach requires a different technique than that used in the present paper, and remains as an open problem.

## 5. Proofs of Theorems 1 and 2.

5.1. *Proof of Theorem 1.* Before proving Theorem 1, we first prepare some technical lemmas (Lemmas 1-3) and establish preliminary rates of contraction for the quasi-posterior distribution (Proposition 4). Proofs of

---

<sup>6</sup>The previous version of this paper contains results on shrinking priors, but  $J_n$  should be still slowly growing as in the above propositions, which corresponds to the sieve method using slowly growing sieves with *light* penalties. Those results are removed in the current version according to the referee’s suggestion, but available upon request.

Lemmas 1-3 are given in Appendix A. For the notational convenience, define the matrices

$$\hat{\Phi}_{WX} = \mathbb{E}_n[\phi^{J_n}(W_i)\phi^{J_n}(X_i)^T], \quad \hat{\Phi}_{XW} = \hat{\Phi}_{WX}^T, \quad \hat{\Phi}_{WW} = \mathbb{E}_n[\phi^{J_n}(W_i)^{\otimes 2}],$$

which are the empirical counterparts of  $\Phi^{WX}$ ,  $\Phi^{XW}$  and  $\Phi_{WW}$ , respectively. Also define

$$U_i = Y_i - g_0(X_i), \quad R_i = Y_i - P_{J_n}g_0(X_i), \quad \Delta_n = \sqrt{n}\mathbb{E}_n[\phi^{J_n}(W_i)R_i].$$

Lemma 1 is a technical lemma on these quantities. Lemma 2 characterizes the total variation convergence between two centered multivariate normal distributions with increasing dimensions in terms of the speed of convergence between the corresponding covariance matrices. Lemma 3 will be used in the latter part in the proof of Theorem 1

LEMMA 1. *Suppose that Assumptions 1-4 are satisfied. Let  $J_n \rightarrow \infty$  as  $n \rightarrow \infty$ . (i) There exists a constant  $D > 0$  such that  $\sup_{w \in [0,1]} \|\phi^J(w)\|_{\ell^2} \leq D2^{J/2}$  for all  $J \geq J_0$ . (ii)  $C_1^{-1} \leq s_{\min}(\mathbb{E}[\phi^J(W)^{\otimes 2}]) \leq s_{\max}(\mathbb{E}[\phi^J(W)^{\otimes 2}]) \leq C_1$  and  $s_{\max}(\mathbb{E}[\phi^J(W)\phi^J(X)^T]) \leq C_1$  for all  $J \geq J_0$ . (iii) If  $J_n2^{J_n}/n \rightarrow 0$ ,  $\|\hat{\Phi}_{WW} - \Phi_{WW}\|_{\text{op}} = O_P(\sqrt{J_n2^{J_n}/n})$  and  $\|\hat{\Phi}_{WX} - \Phi_{WX}\|_{\text{op}} = O_P(\sqrt{J_n2^{J_n}/n})$ . (iv)  $\|\mathbb{E}_n[\phi^{J_n}(W_i)R_i]\|_{\ell^2}^2 = O_P(2^{J_n}/n + \tau_{J_n}^2 2^{-2J_n s})$ . (v) If  $J_n2^{J_n}/n = o(\tau_{J_n}^2)$ ,  $s_{\min}(\hat{\Phi}_{WX}) \geq (1 - o_P(1))\tau_{J_n}$ .*

LEMMA 2. *Let  $\Sigma_n$  be a sequence of symmetric positive definite matrices of dimension  $k_n \rightarrow \infty$  as  $n \rightarrow \infty$  such that  $\|\Sigma_n - I_{k_n}\|_{\text{op}} = o(k_n^{-1})$ . Then as  $n \rightarrow \infty$ ,*

$$\int |dN(0, \Sigma_n)(x) - dN(0, I_{k_n})(x)| dx \rightarrow 0.$$

LEMMA 3. *Let  $\hat{A}_n$  be a sequence of random  $k_n \times k_n$  matrices where  $k_n$  is either bounded or  $k_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Suppose that there exist sequences of positive constants  $\epsilon_n, \delta_n$  and a sequence of non-random, non-singular  $k_n \times k_n$  matrices  $A_n$  such that  $\epsilon_n \rightarrow 0, \delta_n \rightarrow 0, s_{\min}(A_n) \gtrsim \epsilon_n, \|\hat{A}_n - A_n\|_{\text{op}} = O_P(\delta_n)$  and  $\epsilon_n^{-1}\delta_n \rightarrow 0$ . Then  $\hat{A}_n$  is non-singular with probability approaching one and  $\|\hat{A}_n^{-1}A_n - I_{k_n}\|_{\text{op}} \vee \|A_n\hat{A}_n^{-1} - I_{k_n}\|_{\text{op}} = O_P(\epsilon_n^{-1}\delta_n)$ .*

The following proposition gives preliminary rates of contraction for the quasi-posterior distribution.

PROPOSITION 4 (Preliminary contraction rates). *Suppose that Assumptions 1-4 are satisfied. Take  $J_n$  in such a way that  $J_n \rightarrow \infty$  and  $J_n2^{J_n}/n =$*

$o(\tau_{J_n}^2)$ . Let  $\epsilon_n$  be a sequence of positive constants such that  $\epsilon_n \rightarrow 0$  and  $\sqrt{n}\epsilon_n \rightarrow \infty$ . Assume that a sequence of generating priors  $\tilde{\Pi}_n$  satisfies condition P1) of Theorem 1. Define the data-dependent, empirical seminorm  $\|\cdot\|_{\mathcal{D}_n}$  on  $\mathbb{R}^{2^{J_n}}$  by

$$\|b^{J_n}\|_{\mathcal{D}_n} = \|\hat{\Phi}_{WX}b^{J_n}\|_{\ell^2}, \quad b^{J_n} \in \mathbb{R}^{2^{J_n}}.$$

Then for every sequence  $M_n \rightarrow \infty$ , we have

$$\tilde{\Pi}_n\{b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\mathcal{D}_n} > M_n(\epsilon_n + \tau_{J_n}2^{-J_n s}) \mid \mathcal{D}_n\} \xrightarrow{P} 0.$$

PROOF OF PROPOSITION 4. The proof consists of constructing suitable “tests” and is essentially similar to, e.g., the proof of Theorem 2.1 in [24]. Let  $\delta_n = \epsilon_n + \tau_{J_n}2^{-J_n s}$ . We wish to show that there exists a constant  $c_0 > 0$  such that

$$(15) \quad \mathbb{P}\left\{\tilde{\Pi}_n(b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\mathcal{D}_n} > M_n\delta_n \mid \mathcal{D}_n) \leq e^{-c_0 M_n^2 n \delta_n^2}\right\} \rightarrow 1.$$

Note that since  $\sqrt{n}\epsilon_n \rightarrow \infty$ ,  $n\delta_n^2 \geq n\epsilon_n^2 \rightarrow \infty$ . Below,  $c_1, c_2, \dots$  are some positive constants of which the values are understood in the context.

Note that  $Y_i = P_{J_n}g_0(X_i) + R_i = \phi^{J_n}(X_i)^T b_0^{J_n} + R_i$ . Then for  $b^{J_n} \in \mathbb{R}^{2^{J_n}}$ ,

$$(16) \quad \begin{aligned} \mathbb{E}_n[\hat{m}^2(W_i, b^{J_n})] &= -2(b^{J_n} - b_0^{J_n})^T \hat{\Phi}_{XW} \hat{\Phi}_{WW}^{-1} \mathbb{E}_n[\phi^{J_n}(W_i)R_i] \\ &\quad + (b^{J_n} - b_0^{J_n})^T \hat{\Phi}_{XW} \hat{\Phi}_{WW}^{-1} \hat{\Phi}_{WX} (b^{J_n} - b_0^{J_n}) \\ &\quad + \mathbb{E}_n[\phi^{J_n}(W_i)R_i]^T \hat{\Phi}_{WW}^{-1} \mathbb{E}_n[\phi^{J_n}(W_i)R_i]. \end{aligned}$$

Since the last term is independent of  $b^{J_n}$ , it is canceled out in the quasi-posterior distribution. Denote by  $\ell_{b^{J_n}}(\mathcal{D}_n)$  the sum of the first two terms in (16). Then

$$\tilde{\Pi}_n(db^{J_n} \mid \mathcal{D}_n) \propto \exp\{-(n/2)\ell_{b^{J_n}}(\mathcal{D}_n)\} \tilde{\Pi}_n(db^{J_n}).$$

Using the fact that for any  $x, y, c \in \mathbb{R}$  with  $c > 0$ ,  $2xy \leq cx^2 + c^{-1}y^2$ , we have

$$(17) \quad \begin{aligned} \ell_{b^{J_n}}(\mathcal{D}_n) &\geq (\hat{\lambda}_{\min} - c)\|b^{J_n} - b_0^{J_n}\|_{\mathcal{D}_n}^2 \\ &\quad - c^{-1}\hat{\lambda}_{\max}^2\|\mathbb{E}_n[\phi^{J_n}(W_i)R_i]\|_{\ell^2}^2, \quad \forall c > 0, \end{aligned}$$

where  $\hat{\lambda}_{\min}$  and  $\hat{\lambda}_{\max}$  are the minimum and maximum eigenvalues of the matrix  $\hat{\Phi}_{WW}^{-1}$ , respectively. Likewise, we have

$$(18) \quad \begin{aligned} \ell_{b^{J_n}}(\mathcal{D}_n) &\leq (\hat{\lambda}_{\max} + c)\|b^{J_n} - b_0^{J_n}\|_{\mathcal{D}_n}^2 \\ &\quad + c^{-1}\hat{\lambda}_{\max}^2\|\mathbb{E}_n[\phi^{J_n}(W_i)R_i]\|_{\ell^2}^2, \quad \forall c > 0. \end{aligned}$$

Define the event

$$\begin{aligned} \mathcal{E}_{1n} = \{ \mathcal{D}_n : \hat{\lambda}_{\min} < 0.5C_1^{-1} \} \cup \{ \mathcal{D}_n : \hat{\lambda}_{\max} > 1.5C_1 \} \\ \cup \{ \mathcal{D}_n : \|\mathbb{E}_n[\phi^{J_n}(W_i)R_i]\|_{\ell^2}^2 > M_n\delta_n^2 \}. \end{aligned}$$

Construct the “tests”  $\omega_n$  by  $\omega_n = 1(\mathcal{E}_{1n})$ . Then we have

$$\begin{aligned} & \tilde{\Pi}_n(b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\mathcal{D}_n} > M_n\delta_n \mid \mathcal{D}_n) \\ & = \tilde{\Pi}_n(b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\mathcal{D}_n} > M_n\delta_n \mid \mathcal{D}_n) \{ \omega_n + (1 - \omega_n) \} \\ (19) \quad & \leq \omega_n + \tilde{\Pi}_n(b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\mathcal{D}_n} > M_n\delta_n \mid \mathcal{D}_n) (1 - \omega_n). \end{aligned}$$

By Lemmas 1 (ii)-(iv), we have  $\mathbb{P}(\omega_n = 1) = \mathbb{P}(\mathcal{E}_{1n}) \rightarrow 0$ .

For the second term in (19), taking  $c > 0$  sufficiently small in (17), we have

$$\begin{aligned} & (1 - \omega_n) \int_{\|b^{J_n} - b_0^{J_n}\|_{\mathcal{D}_n} > M_n\delta_n} \exp\{-(n/2)\ell_{b^{J_n}}(\mathcal{D}_n)\} \tilde{\Pi}_n(db^{J_n}) \\ & \leq \exp\{-c_1M_n^2n\delta_n^2 + O(M_nn\delta_n^2)\} \leq e^{-c_2M_n^2n\delta_n^2}. \end{aligned}$$

On the other hand, taking, say  $c = 1$  in (18), we have

$$\begin{aligned} & (1 - \omega_n) \int \exp\{-(n/2)\ell_{b^{J_n}}(\mathcal{D}_n)\} \tilde{\Pi}_n(db^{J_n}) \\ & \geq (1 - \omega_n) \int_{\|b^{J_n} - b_0^{J_n}\|_{\mathcal{D}_n} \leq \sqrt{M_n}\epsilon_n} \exp\{-(n/2)\ell_{b^{J_n}}(\mathcal{D}_n)\} \tilde{\Pi}_n(db^{J_n}) \\ & \geq (1 - \omega_n) e^{-c_3M_nn\delta_n^2} \int_{\|b^{J_n} - b_0^{J_n}\|_{\mathcal{D}_n} \leq \sqrt{M_n}\epsilon_n} \tilde{\Pi}_n(db^{J_n}). \end{aligned}$$

Denote by  $\hat{s}_{\max}$  the maximum singular value of the matrix  $\hat{\Phi}_{WX}$ , so that

$$\|b^{J_n} - b_0^{J_n}\|_{\mathcal{D}_n} \leq \hat{s}_{\max} \|b^{J_n} - b_0^{J_n}\|_{\ell^2}.$$

Define the event  $\mathcal{E}_{2n} = \{ \mathcal{D}_n : \hat{s}_{\max} \leq 1.5C_1 \}$ . By Lemmas 1 (ii) and (iii), we have  $\mathbb{P}(\mathcal{E}_{2n}) \rightarrow 1$ . Since  $M_n \rightarrow \infty$ , for all  $n$  sufficiently large, we have

$$\begin{aligned} & 1(\mathcal{E}_{2n})(1 - \omega_n) \int \exp\{-(n/2)\ell_{b^{J_n}}(\mathcal{D}_n)\} \tilde{\Pi}_n(db^{J_n}) \\ & \geq 1(\mathcal{E}_{2n})(1 - \omega_n) e^{-c_3M_nn\delta_n^2} \tilde{\Pi}_n(b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\ell^2} \leq \epsilon_n) \\ & \geq 1(\mathcal{E}_{2n})(1 - \omega_n) e^{-c_3M_n\delta_n^2 - Cn\epsilon_n^2} \\ & \geq 1(\mathcal{E}_{2n})(1 - \omega_n) e^{-c_4M_nn\delta_n^2}, \end{aligned}$$

where the second inequality is due to the small ball condition P1). Summarizing, we have

$$\tilde{\Pi}_n(b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\mathcal{D}_n} > M_n \delta_n \mid \mathcal{D}_n)(1 - \omega_n) \leq 1(\mathcal{E}_{2n}^c) + e^{-c_2 M_n^2 n \delta_n^2 + c_4 M_n n \delta_n^2}.$$

Therefore, we obtain (15) for a sufficiently small  $c_0 > 0$ .  $\square$

We are now in position to prove Theorem 1. We will say that a sequence of random variables  $A_n$  is *eventually* bounded by another sequence of random variables  $B_n$  if  $\mathbb{P}(A_n \leq B_n) \rightarrow 1$  as  $n \rightarrow \infty$ .

PROOF OF THEOREM 1. We first note that by Lemmas 1 (ii), (iii) and (v), the matrices  $\hat{\Phi}_{WX}$  and  $\hat{\Phi}_{WW}$  are non-singular with probability approaching one. Conditional on  $\mathcal{D}_n$ , define the rescaled ‘‘parameter’’  $\theta^{J_n} = (\theta_1, \dots, \theta_{2^{J_n}})^T = \sqrt{n} \hat{\Phi}_{WX}(b^{J_n} - b_0^{J_n})$ . By (16), the corresponding ‘‘quasi-posterior’’ density for  $\theta^{J_n}$  is given by

$$\pi_n^*(\theta^{J_n} \mid \mathcal{D}_n) d\theta^{J_n} \propto \tilde{\pi}_n(b_0^{J_n} + \hat{\Phi}_{WX}^{-1} \theta^{J_n} / \sqrt{n}) dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n}) d\theta^{J_n},$$

where recall that  $\Delta_n = \sqrt{n} \mathbb{E}_n[\phi^{J_n}(W_i) R_i]$  (this operation is valid as soon as  $\hat{\Phi}_{WX}$  and  $\hat{\Phi}_{WW}$  are non-singular, of which the probability is approaching one).

The proof of Theorem 1 consists of 3 steps. After Step 1, we will turn to the proof of (9). The remaining two steps are devoted to the proof of (10).

Step 1. We first show that

$$(20) \quad \int |\pi_n^*(\theta^{J_n} \mid \mathcal{D}_n) - dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})| d\theta^{J_n} \xrightarrow{P} 0.$$

In this step, we do *not* assume  $J_n 2^{3J_n} / n = o(\tau_{J_n}^2)$ . As before, let  $\delta_n = \epsilon_n + \tau_{J_n} 2^{-J_n^s}$ . By Proposition 4, for every sequence  $M_n \rightarrow \infty$ ,

$$\int_{\|\theta^{J_n}\|_{\ell^2} \leq M_n \sqrt{n} \delta_n} \pi_n^*(\theta^{J_n} \mid \mathcal{D}_n) d\theta^{J_n} = 1 + o_P(1),$$

by which we have

$$(21) \quad \begin{aligned} & \text{Left side of (20)} \\ & \leq \int_{\|\theta^{J_n}\|_{\ell^2} \leq M_n \sqrt{n} \delta_n} |\pi_n^*(\theta^{J_n} \mid \mathcal{D}_n) - dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})| d\theta^{J_n} \\ & \quad + \int_{\|\theta^{J_n}\|_{\ell^2} > M_n \sqrt{n} \delta_n} dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n}) d\theta^{J_n} + o_P(1). \end{aligned}$$

By Lemma 1 (iv),  $\|\Delta_n\|_{\ell^2} = O_P(\sqrt{n}\delta_n)$ , and by Lemmas 1 (ii) and (iii),  $(1 - o_P(1))C_1^{-1} \leq s_{\min}(\hat{\Phi}_{WW}) \leq s_{\max}(\hat{\Phi}_{WW}) \leq (1 + o_P(1))C_1$ , so that the second integral is eventually bounded by

$$(22) \quad \int_{\|\theta^{J_n}\|_{\ell^2} > \sqrt{M_n n} \delta_n} dN(0, I_{2^{J_n}})(\theta^{J_n}) d\theta^{J_n},$$

where note that  $M_n$  is replaced by  $\sqrt{M_n}$  to “absorb” the constant. By Borell’s inequality for Gaussian measures [see, for example, 60, Lemma A.2.2], for every  $x > 0$ ,

$$(23) \quad \mathbb{P}(\|N(0, I_{2^{J_n}})\|_{\ell^2} > \sqrt{2^{J_n}} + x) \leq e^{-x^2/2}.$$

Here since  $n\delta_n^2 \geq n\epsilon_n^2 \gtrsim 2^{J_n}$ ,  $\sqrt{M_n n} \delta_n / \sqrt{2^{J_n}} \rightarrow \infty$ , so that the integral in (22) is  $o(1)$ .

It remains to show that the first integral in (21) is  $o_P(1)$ . This step uses a standard cancellation argument. Let  $\mathcal{C}_n := \{\theta^{J_n} \in \mathbb{R}^{2^{J_n}} : \|\theta^{J_n}\|_{\ell^2} \leq M_n \sqrt{n} \delta_n\}$ . First, provided that  $\|\hat{\Phi}_{WX}^{-1}\|_{\text{op}} \leq 1.5\tau_{J_n}^{-1}$ , for all  $\theta^{J_n} \in \mathcal{C}_n$ ,

$$\|\hat{\Phi}_{WX}^{-1} \theta^{J_n} / \sqrt{n}\|_{\ell^2} \leq 1.5M_n \tau_{J_n}^{-1} \delta_n \leq 1.5M_n (2^{-J_n s} + \tau_{J_n}^{-1} \epsilon_n) \sim M_n \gamma_n.$$

So taking  $M_n \rightarrow \infty$  such that  $M_n = o(L_n)$ ,  $\|\hat{\Phi}_{WX}^{-1} \theta^{J_n} / \sqrt{n}\|_{\ell^2} \leq L_n \gamma_n$  and hence  $\tilde{\pi}_n(b_0^{J_n} + \hat{\Phi}_{WX}^{-1} \theta^{J_n} / \sqrt{n}) > 0$  for all  $n$  sufficiently large. Here, by Lemma 1 (v), we have  $\mathbb{P}(\|\hat{\Phi}_{WX}^{-1}\|_{\text{op}} \leq 1.5\tau_{J_n}^{-1}) \rightarrow 1$ .

Suppose that  $\|\hat{\Phi}_{WX}^{-1}\|_{\text{op}} \leq 1.5\tau_{J_n}^{-1}$ . Let

$$\pi_{n, \mathcal{C}_n}^*(\theta^{J_n} \mid \mathcal{D}_n) \text{ and } dN^{\mathcal{C}_n}(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})$$

denote the probability densities obtained by first restricting  $\pi_n^*(\theta^{J_n} \mid \mathcal{D}_n)$  and  $dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})$  to the ball  $\mathcal{C}_n$  and then renormalizing, respectively. By the first part of the present proof, replacing  $\pi_n^*(\theta^{J_n} \mid \mathcal{D}_n)$  and  $dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})$  by  $\pi_{n, \mathcal{C}_n}^*(\theta^{J_n} \mid \mathcal{D}_n)$  and  $dN^{\mathcal{C}_n}(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})$  respectively in the first integral in (21) has impact at most  $o_P(1)$ . Abbreviating  $\pi_{n, \mathcal{C}_n}^*(\theta^{J_n} \mid \mathcal{D}_n)$  by  $\pi_{n, \mathcal{C}_n}^*$ ,  $dN^{\mathcal{C}_n}(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})$  by  $dN^{\mathcal{C}_n}$ ,  $dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})$  by  $dN$ , and  $\tilde{\pi}_n(b_0^{J_n} + \hat{\Phi}_{WX}^{-1} \theta^{J_n} / \sqrt{n})$  by  $\tilde{\pi}_n$ , we have

$$\begin{aligned} \int |\pi_{n, \mathcal{C}_n}^* - dN^{\mathcal{C}_n}| &= \int \left| 1 - \frac{dN^{\mathcal{C}_n}}{\pi_{n, \mathcal{C}_n}^*} \right| \pi_{n, \mathcal{C}_n}^* = \int \left| 1 - \frac{dN / \int_{\mathcal{C}_n} dN}{\tilde{\pi}_n dN / \int_{\mathcal{C}_n} \tilde{\pi}_n dN} \right| \pi_{n, \mathcal{C}_n}^* \\ &= \int \left| 1 - \frac{\int_{\mathcal{C}_n} \tilde{\pi}_n dN}{\tilde{\pi}_n \int_{\mathcal{C}_n} dN} \right| \pi_{n, \mathcal{C}_n}^* = \int \left| 1 - \frac{\int_{\mathcal{C}_n} \tilde{\pi}_n dN^{\mathcal{C}_n}}{\tilde{\pi}_n} \right| \pi_{n, \mathcal{C}_n}^*. \end{aligned}$$

By the convexity of the map  $x \mapsto |1 - x|$  and Jensen's inequality, the last expression is bounded by

$$\sup_{\theta^{J_n} \in \mathcal{C}_n, \tilde{\theta}^{J_n} \in \mathcal{C}_n} \left| 1 - \frac{\tilde{\pi}_n(b_0^{J_n} + \hat{\Phi}_{WX}^{-1} \theta^{J_n} / \sqrt{n})}{\tilde{\pi}_n(b_0^{J_n} + \hat{\Phi}_{WX}^{-1} \tilde{\theta}^{J_n} / \sqrt{n})} \right|,$$

which is eventually bounded by

$$\sup_{\|b^{J_n}\|_{\ell^2} \leq L_n \gamma_n, \|\tilde{b}^{J_n}\|_{\ell^2} \leq L_n \gamma_n} \left| 1 - \frac{\tilde{\pi}_n(b_0^{J_n} + b^{J_n})}{\tilde{\pi}_n(b_0^{J_n} + \tilde{b}^{J_n})} \right|.$$

The last expression goes to zeros as  $n \rightarrow \infty$  by condition P2).

We now turn to the proof of (9). Take any  $M_n \rightarrow \infty$  (this  $M_n$  may be different from the previous  $M_n$ ). By Step 1, we have

$$\begin{aligned} & \sup_{z > 0} \left| \tilde{\Pi}_n \{ b^{J_n} : \|\hat{\Phi}_{WX}(b^{J_n} - b_0^{J_n})\|_{\ell^2} > z \mid \mathcal{D}_n \} \right. \\ & \quad \left. - \int_{\|\theta^{J_n}\|_{\ell^2} > z} dN(n^{-1/2} \Delta_n, n^{-1} \hat{\Phi}_{WW})(\theta^{J_n}) d\theta^{J_n} \right| \xrightarrow{P} 0. \end{aligned}$$

By Lemma 1 (v), we have

$$\begin{aligned} \|\hat{\Phi}_{WX}(b^{J_n} - b_0^{J_n})\|_{\ell^2} & \geq s_{\min}(\hat{\Phi}_{WX}) \|b^{J_n} - b_0^{J_n}\|_{\ell^2} \\ & \geq (1 - o_P(1)) \tau_{J_n} \|b^{J_n} - b_0^{J_n}\|_{\ell^2}, \end{aligned}$$

by which we have, uniformly in  $z > 0$ ,

$$\begin{aligned} & \tilde{\Pi}_n \{ b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\ell^2} > 2\tau_{J_n}^{-1} z \mid \mathcal{D}_n \} \\ & \leq \tilde{\Pi}_n \{ b^{J_n} : \|\hat{\Phi}_{WX}(b^{J_n} - b_0^{J_n})\|_{\ell^2} > z \mid \mathcal{D}_n \} + o_P(1) \\ & \leq \int_{\|\theta^{J_n}\|_{\ell^2} > z} dN(n^{-1/2} \Delta_n, n^{-1} \hat{\Phi}_{WW})(\theta^{J_n}) d\theta^{J_n} + o_P(1). \end{aligned}$$

By Markov's inequality, the integral in the last expression is bounded by

$$\frac{1}{nz^2} \{ \|\Delta_n\|_{\ell^2}^2 + \text{tr}(\hat{\Phi}_{WW}) \}.$$

By Lemmas 1 (ii)-(iv), we have  $\|\Delta_n\|_{\ell^2}^2 + \text{tr}(\hat{\Phi}_{WW}) = O_P(2^{J_n} + n\tau_{J_n}^2 2^{-2J_n s})$ . Therefore, we conclude that, taking  $z = M_n(\tau_{J_n} 2^{-J_n s} + \sqrt{2^{J_n}/n})$ ,  $\tilde{\Pi}_n \{ b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\ell^2} > 2M_n(2^{-J_n s} + \tau_{J_n}^{-1} \sqrt{2^{J_n}/n}) \mid \mathcal{D}_n \} \xrightarrow{P} 0$ , which leads to the contraction rate result (9).

In what follows, we assume  $J_n 2^{3J_n}/n = o(\tau_{J_n}^2)$ , and prove the asymptotic normality result (10).

Step 2. (Replacement of  $\hat{\Phi}_{WW}$  by  $\Phi_{WW}$ ). This step shows that

$$\int |dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n}) - dN(\Delta_n, \Phi_{WW})(\theta^{J_n})| d\theta^{J_n} \xrightarrow{P} 0,$$

which is equivalent to

$$\int |dN(0, \hat{\Phi}_{WW})(\theta^{J_n}) - dN(0, \Phi_{WW})(\theta^{J_n})| d\theta^{J_n} \xrightarrow{P} 0.$$

By Lemmas 1 (ii), (iii) and Lemma 2, this follows if  $\sqrt{J_n 2^{J_n}/n} = o(2^{-J_n})$ , i.e.,  $J_n 2^{3J_n} = o(n)$ , which is satisfied since  $J_n 2^{3J_n}/n = o(\tau_{J_n}^2) = o(1)$ .

Step 3. (Replacement of  $\hat{\Phi}_{WX}$  by  $\Phi_{WX}$ ). We have shown that

$$\int |\pi_n^*(\theta^{J_n} | \mathcal{D}_n) - dN(\Delta_n, \Phi_{WW})(\theta^{J_n})| d\theta^{J_n} \xrightarrow{P} 0.$$

By Scheffé's lemma, this means that

$$\|\tilde{\Pi}_n\{b^{J_n} : \sqrt{n}\hat{\Phi}_{WX}(b^{J_n} - b_0^{J_n}) \in \cdot | \mathcal{D}_n\} - N(\Delta_n, \Phi_{WW})(\cdot)\|_{\text{TV}} \xrightarrow{P} 0,$$

or equivalently,

$$\|\tilde{\Pi}_n\{b^{J_n} : \sqrt{n}(b^{J_n} - b_0^{J_n}) \in \cdot | \mathcal{D}_n\} - N(\hat{\Phi}_{WX}^{-1}\Delta_n, \hat{\Phi}_{WX}^{-1}\Phi_{WW}\hat{\Phi}_{XW}^{-1})(\cdot)\|_{\text{TV}} \xrightarrow{P} 0.$$

The last expression is asymptotically valid since  $\hat{\Phi}_{WX}$  is non-singular with probability approaching one. Recall the maximum quasi-likelihood estimator  $\hat{b}^{J_n}$ . With probability approaching one, we have

$$\hat{b}^{J_n} = \hat{\Phi}_{WX}^{-1} \mathbb{E}_n[\phi^{J_n}(W_i)Y_i] = b_0^{J_n} + \hat{\Phi}_{WX}^{-1} \mathbb{E}_n[\phi^{J_n}(W_i)R_i],$$

so that  $\sqrt{n}(\hat{b}^{J_n} - b_0^{J_n}) = \hat{\Phi}_{WX}^{-1}\Delta_n$ . Hence to conclude the theorem, it suffices to show that

$$(24) \quad \|N(\hat{\Phi}_{WX}^{-1}\Delta_n, \hat{\Phi}_{WX}^{-1}\Phi_{WW}\hat{\Phi}_{XW}^{-1}) - N(\hat{\Phi}_{WX}^{-1}\Delta_n, \hat{\Phi}_{WX}^{-1}\Phi_{WW}\hat{\Phi}_{XW}^{-1})\|_{\text{TV}} \xrightarrow{P} 0.$$

Assertion (24) reduces to

$$\|N(0, \hat{\Phi}_{WX}^{-1}\Phi_{WW}\hat{\Phi}_{XW}^{-1}) - N(0, \Phi_{WW})\|_{\text{TV}} \xrightarrow{P} 0.$$

By Lemmas 1 (ii), (iii) and Lemma 3,  $\|\hat{\Phi}_{WX}^{-1}\Phi_{WW}\hat{\Phi}_{XW}^{-1} - \Phi_{WW}\|_{\text{op}} = O_P(\tau_{J_n}^{-1} \sqrt{J_n 2^{J_n}/n}) = o_P(2^{-J_n})$  (the last equality follows since  $J_n 2^{3J_n}/n = o(\tau_{J_n}^2)$ ). Since  $C_1^{-1} \leq s_{\min}(\Phi_{WW}) \leq s_{\max}(\Phi_{WW}) \leq C_1$ , the desired conclusion follows from Lemma 2.

Steps 1-3 lead to the asymptotic normality result (10).  $\square$

5.2. *Proof of Theorem 2.* We first prove the following lemma.

LEMMA 4. *Suppose that the conditions of Theorem 2 are satisfied. Then there exists a constant  $D > 0$  such that*

$$\mathbb{P} \left\{ \|\mathbb{E}_n[\phi^{J_n}(W_i)U_i]\|_{\ell^2} > D\sqrt{2^{J_n}/n} \right\} \rightarrow 0.$$

REMARK 9. It is standard to show that  $\|\mathbb{E}_n[\phi^{J_n}(W_i)U_i]\|_{\ell^2} = O_P(\sqrt{2^{J_n}/n})$ , which, however, does not lead to the conclusion of Lemma 4 since the former only implies that for every sequence  $M_n \rightarrow \infty$ ,  $\mathbb{P}\{\|\mathbb{E}_n[\phi^{J_n}(W_i)U_i]\|_{\ell^2} > M_n\sqrt{2^{J_n}/n}\} \rightarrow 0$ . Hence an additional step is needed. The current proof uses a truncation argument and Talagrand's concentration inequality.

PROOF OF LEMMA 4. For a given  $\lambda > 0$ , define  $U_i^- = U_i 1(|U_i| \leq \lambda)$  and  $U_i^+ = U_i 1(|U_i| > \lambda)$ . Since  $0 = \mathbb{E}[U | W] = \mathbb{E}[U^- | W] + \mathbb{E}[U^+ | W]$ , we have  $\mathbb{E}_n[\phi^{J_n}(W_i)U_i] = n^{-1} \sum_{i=1}^n \{\phi^{J_n}(W_i)U_i^- - \mathbb{E}[\phi^{J_n}(W)U^-]\} + n^{-1} \sum_{i=1}^n \{\phi^{J_n}(W_i)U_i^+ - \mathbb{E}[\phi^{J_n}(W)U^+]\}$ , by which we have

$$\begin{aligned} \|\mathbb{E}_n[\phi^{J_n}(W_i)U_i]\|_{\ell^2} &\leq \|n^{-1} \sum_{i=1}^n \{\phi^{J_n}(W_i)U_i^- - \mathbb{E}[\phi^{J_n}(W)U^-]\}\|_{\ell^2} \\ &\quad + \|n^{-1} \sum_{i=1}^n \{\phi^{J_n}(W_i)U_i^+ - \mathbb{E}[\phi^{J_n}(W)U^+]\}\|_{\ell^2} \\ &=: I + II. \end{aligned}$$

First, by Markov's inequality, we have for every  $z > 0$ ,

$$\begin{aligned} \mathbb{P}(II > z) &\leq \frac{\mathbb{E}[II^2]}{z^2} \leq \frac{\sum_{l=1}^{2^{J_n}} \mathbb{E}[(\phi_l(W)U^+)^2]}{nz^2} \\ &\leq \frac{\sup_{w \in [0,1]} \mathbb{E}[U^2 1(|U| > \lambda) | W = w] \times \sum_{l=1}^{2^{J_n}} \mathbb{E}[\phi_l(W)^2]}{nz^2} \\ &\leq \frac{C_1 2^{J_n}}{nz^2} \times \sup_{w \in [0,1]} \mathbb{E}[U^2 1(|U| > \lambda) | W = w], \end{aligned}$$

where we have used that  $\sum_{l=1}^{2^{J_n}} \mathbb{E}[\phi_l(W)^2] = \text{tr}(\Phi_{WW}) \leq 2^{J_n} s_{\max}(\Phi_{WW}) \leq C_1 2^{J_n}$  by Lemma 1 (ii). Thus we have

$$\mathbb{P}\{II > \sqrt{C_1 2^{J_n}/n}\} \leq \sup_{w \in [0,1]} \mathbb{E}[U^2 1(|U| > \lambda) | W = w].$$

By assumption, the right side goes to zero as  $\lambda \rightarrow \infty$ .

Second, let  $Z_i = \phi^{J_n}(W_i)U_i^- - \mathbb{E}[\phi^{J_n}(W)U^-]$  (denote by  $Z$  the generic version of  $Z_i$ ). Let  $\mathbb{S}^{2^{J_n}-1} := \{\alpha^{J_n} \in \mathbb{R}^{2^{J_n}} : \|\alpha^{J_n}\|_{\ell^2} = 1\}$ . Then

$$I = \|\mathbb{E}_n[Z_i]\|_{\ell^2} = \sup_{\alpha^{J_n} \in \mathbb{S}^{2^{J_n}-1}} \mathbb{E}_n[(\alpha^{J_n})^T Z_i].$$

We make use of Talagrand's concentration inequality to bound the tail probability of  $I$ . For any  $\alpha^{J_n} \in \mathbb{S}^{2^{J_n}-1}$ , by Lemma 1, we have

$$\begin{aligned} \mathbb{E}\{[(\alpha^{J_n})^T Z]^2\} &\leq \sup_{w \in [0,1]} \mathbb{E}[U^2 \mid W = w] \times s_{\max}(\Phi_{WW}) \leq C_1^2, \\ |(\alpha^{J_n})^T Z| &\leq \lambda \sup_{w \in [0,1]} \|\phi^{J_n}(w)\|_{\ell^2} \leq D_1 \lambda \sqrt{2^{J_n}}, \text{ and} \\ (\mathbb{E}[I])^2 \leq \mathbb{E}[I^2] &\leq n^{-1} \sup_{w \in [0,1]} \mathbb{E}[U^2 \mid W = w] \times \sum_{l=1}^{2^{J_n}} \mathbb{E}[\phi_l(W)^2] \leq C_1^2 2^{J_n}/n, \end{aligned}$$

where  $D_1 > 0$  is a constant. Thus, by Talagrand's inequality (see Theorem 5 in Appendix E), we have for every  $z > 0$ ,

$$\mathbb{P}\{I \geq D_2(\sqrt{2^{J_n}/n} + \sqrt{z/n} + z\lambda\sqrt{2^{J_n}/n})\} \leq e^{-z},$$

where  $D_2 > 0$  is a constant independent of  $\lambda$  and  $z$ .

The final conclusion follows from taking  $\lambda = \lambda_n \rightarrow \infty$  and  $z = z_n \rightarrow \infty$  sufficiently slowly.  $\square$

PROOF OF THEOREM 2. Let  $D_1$  and  $D_2$  be some positive constants of which the values are understood in the context. For either  $g_0 \in B_{\infty,\infty}^s$  or  $B_{2,2}^s$ ,  $\|g_0 - P_{J_n}g_0\| = O(2^{-J_n s}) = o(1)$ , by which we have

$$\begin{aligned} \sum_{l=1}^{2^{J_n}} \text{Var}\{\mathbb{E}_n[\phi_l(W_i)(g_0 - P_{J_n}g_0)(X_i)]\} &\leq n^{-1} \sum_{l=1}^{2^{J_n}} \mathbb{E}[\phi_l(W)^2 \{(g_0 - P_{J_n}g_0)(X)\}^2] \\ &= n^{-1} \sum_{l=1}^{2^{J_n}} \iint \phi_l(w)^2 \{(g_0 - P_{J_n}g_0)(x)\}^2 f_{X,W}(x,w) dx dw \\ &\leq n^{-1} C_1 \|g_0 - P_{J_n}g_0\|^2 \times \sum_{l=1}^{2^{J_n}} \int \phi_l(w)^2 dw = o(2^{J_n}/n). \end{aligned}$$

Hence

$$\begin{aligned} \mathbb{E}_n[\phi^{J_n}(W_i)R_i] &= \mathbb{E}_n[\phi^{J_n}(W_i)U_i] \\ &\quad + \mathbb{E}[\phi^{J_n}(W)(g_0 - P_n g_0)(X)] + \text{Rem}, \end{aligned}$$

with  $\|\text{Rem}\|_{\ell^2} = o_P(\sqrt{2^{J_n}/n})$ . The second term on the right side is  $O(\tau_{J_n} 2^{-J_n s})$  in the Euclidean norm. Together with Lemma 4, we have

$$\mathbb{P}\{\|\mathbb{E}_n[\phi^{J_n}(W_i)R_i]\|_{\ell^2}^2 > D_1(\tau_{J_n}^2 2^{-2J_n s} + 2^{J_n}/n)\} \rightarrow 0.$$

Moreover, by Lemma 1, we have

$$\mathrm{tr}(\hat{\Phi}_{WW}) \leq 2^{J_n} s_{\max}(\hat{\Phi}_{WW}) \leq C_1(1 + o_P(1))2^{J_n}.$$

Taking these together, we have

$$\mathbb{P} \left\{ \|\mathbb{E}_n[\phi^{J_n}(W_i)R_i]\|_{\ell^2}^2 + n^{-1} \mathrm{tr}(\hat{\Phi}_{WW}) \leq D_2(\tau_{J_n}^2 2^{-2J_n s} + 2^{J_n}/n) \right\} \rightarrow 1.$$

By the proof of Theorem 1, this leads to the desired conclusion.  $\square$

**6. Discussion.** We have studied the asymptotic properties of quasi-posterior distributions against sieve priors in the NPIV model and given some specific priors for which the quasi-posterior distribution (the quasi-Bayes estimator) attains the minimax optimal rate of contraction (convergence, respectively). These results greatly sharpen the previous work [46]. We end this paper with two additional discussions.

*6.1. Multivariate case.* In this paper, we have focused on the case where  $X$  and  $W$  are scalar, mainly to avoid the notational complication. It is not difficult to see that the results naturally extend to the case where  $X$  and  $W$  are vectors with the same dimension, by considering tensor product sieves (the contraction/convergence rates will then deteriorate as the dimension grows). We can also consider the following more general situation as in Section 3 of [28]: suppose that  $Y$  is a scalar random variable,  $X$  and  $W$  are random vectors with the same dimension, and  $Z$  is another random vector (whose dimension may be different from  $X$ ), and suppose that we are interested in estimating the function  $g_0$  identified by the conditional moment restriction:  $\mathbb{E}[Y | Z, W] = \mathbb{E}[g_0(X, Z) | Z, W]$  or  $Y = g_0(X, Z) + U$  with  $\mathbb{E}[U | Z, W] = 0$  (i.e.,  $X$  and  $Z$  are endogenous and exogenous explanatory variables, respectively). In principle, the analysis can be reduced to the case where there are no exogenous variables by conditioning on  $Z = z$  (so the sieve measure of ill-posedness can be defined by the one conditional on  $Z = z$ ). More precisely, when  $Z$  is discretely distributed with finitely many mass points, then  $g_0(x, z)$ , where  $z$  is a mass point, can be estimated by using only observations  $i$  for which  $Z_i = z$ . When  $Z$  is continuously distributed, then  $g_0(x, z)$  can be estimated by using observations  $i$  for which  $Z_i$  is “close” to  $z$ ; one way is to use kernel weights as in Section 4.2 of [33]. However, the detailed analysis of this case is not presented here for brevity.

*6.2. Direction of future research.* Lastly we make some remarks on the direction of future research. First, as also noted by [46], (adaptive) selection

of the resolution level  $J_n$  in a (quasi-)Bayesian or “empirical” Bayesian approach is an important topic to be investigated. Second, a (quasi-)Bayesian analysis is typically useful in the analysis of complex models in which frequentist estimation is difficult to implement due to non-differentiability/non-convex nature of loss functions. This usefulness comes from the fact that a (quasi-)Bayesian approach is typically able to avoid numerical optimization. See [13] and [47] for the finite dimensional case. In infinite dimensional models, such a computational challenge in frequentist estimation occurs in the analysis of nonparametric instrumental quantile regression models [35, 11, 21]. In that model, a typical loss function contains the indicator function and hence highly non-convex. In such a case, the computation of an optimal solution is by itself difficult, and a solution obtained, if possible, is typically not guaranteed to be globally optimal since there may be many local optima. It is hence of interest to extend the results of the paper to nonparametric instrumental quantile regression models. The extension to the quantile regression case, which is currently under investigation, is highly non-trivial since the problem of estimating the structural function becomes a *non-linear* ill-posed inverse problem and a delicate care of the stochastic expansion of the criterion function is needed.

**Acknowledgments.** A major part of the work was done while the author was visiting the department of economics, MIT. He would like to thank Professor Victor Chernozhukov for his suggestions and encouragements, as well as Professor Yukitoshi Matsushita for his constructive comments. Also he would like to thank the editor Professor Runze Li, the AE, and anonymous referees for their insightful comments that helped improve on the quality of the paper.

## References.

- [1] Agapiou, S., Larsson, S., and Stuart, A.M. (2013). Posterior contraction rates for the Bayesian approach to linear ill-posed inverse problems. *Stochastic Process. Appl.*, to appear.
- [2] Ai, C., and Chen, X. (2003). Efficient estimation of models with conditional moment restrictions containing unknown functions. *Econometrica* **71** 1795-1843.
- [3] Belloni, A. and Chernozhukov, V. (2009a). On the computational complexity of MCMC-based estimators in large samples. *Ann. Statist.* **37** 2011-2055.
- [4] Belloni, A. and Chernozhukov, V. (2009b). Posterior inference in curved exponential families under increasing dimensions. arXiv:0904.3132. Preprint.
- [5] Bhatia, R. (1997). *Matrix Analysis*. Springer.
- [6] Blundell, R., Chen, X. and Kristensen, D. (2007). Semi-nonparametric IV estimation of shape-invariant Engel curves. *Econometrica* **75** 1613-1669.
- [7] Bontemps, D. (2011). Bernstein-von Mises theorems for Gaussian regression with increasing number of regressors. *Ann. Statist.* **39** 2557-2584.

- [8] Boucheron, S. and Gassiat, E. (2009). A Bernstein-von Mises Theorem for discrete probability distributions. *Electron. J. Statist.* **3** 114-148.
- [9] Carrasco, M., Florens, J.-P., and Renault, E. (2007). Linear inverse problems in structural econometrics: estimation based on spectral decomposition and regularization. In: *Handbook of Econometrics Vol.6* (eds. J.J. Heckman and E.E. Leamer) Elsevier pp. 5633-5751.
- [10] Cavalier, L. (2008). Nonparametric statistical inverse problems. *Inverse Problems* **24** 1-19.
- [11] Chen, X. and Pouzo, D. (2012). Estimation of nonparametric conditional moment models with possibly nonsmooth generalized residuals. *Econometrica* **80** 277-321.
- [12] Chen, X. and Reiss, M. (2011) On rate optimality for ill-posed inverse problems in econometrics. *Econometric Theory* **27** 497-521.
- [13] Chernozhukov, V. and Hong, H. (2003). An MCMC approach to classical estimation. *J. Econometrics* **115** 293-346.
- [14] Clarke, B.S. and Ghosal, S. (2010). Reference priors for exponential families with increasing dimension. *Electron. J. Statist.* **4** 737-780.
- [15] Cohen, A., Daubechies, I., and Vial, P. (1993). Wavelets on the interval and fast wavelet transforms. *Appl. Comput. Harmon. Anal.* **1** 54-81.
- [16] Cox, D. (1993). An analysis of Bayesian inference for nonparametric regression. *Ann. Statist.* **21** 903-923.
- [17] Darolles, S., Fan, Y., Florens, J.P. and Renault, E. (2011). Nonparametric instrumental regression. *Econometrica* **79** 1541-1565.
- [18] Florens, J.-P. and Simoni, A. (2010). Regularized priors for linear inverse problems. IDEI Working Paper n.767-2013.
- [19] Florens, J.-P. and Simoni, A. (2012a). Regularized posteriors in linear ill-posed inverse problems. *Scand. J. Statist.* **39** 214-235.
- [20] Florens, J.-P. and Simoni, A. (2012b). Nonparametric estimation of an instrumental variables regression: a quasi-Bayesian approach based on regularized posterior. *J. Econometrics* **170** 458-475.
- [21] Gagliardini, P. and Scaillet, O. (2011). Nonparametric instrumental variables estimation of structural quantile effects. *Econometrica* **80** 1533-1562.
- [22] Ghosal, S. (1999). Asymptotic normality of posterior distributions in high dimensional linear models. *Bernoulli* **5** 315-331.
- [23] Ghosal, S. (2000). Asymptotic normality of posterior distributions for exponential families when the number of parameters tends to infinity. *J. Multivariate Anal.* **73** 49-68.
- [24] Ghosal, S., Ghosh, J. K. and van der Vaart, A.W. (2000). Convergence rates of posterior distributions. *Ann. Statist.* **28** 500-531.
- [25] Ghosal, S. and van der Vaart, A.W. (2007). Convergence rates of posterior distributions for non-i.i.d. observations. *Ann. Statist.* **35** 192-223.
- [26] Ghosh, J. and Ramamoorthi, R. (2003). *Bayesian Nonparametrics*. Springer.
- [27] Giné, E. and Nickl, R. (2011). Rates of contraction for posterior distributions in  $L^R$  metrics,  $1 \leq R \leq \infty$ . *Ann. Statist.* **39** 2883-2911.
- [28] Hall, P. and Horowitz, J.L. (2005). Nonparametric methods for inference in the presence of instrumental variables. *Ann. Statist.* **33** 2904-2929.
- [29] Härdle, W., Kerkycharian, F., Picard, D., and Tsybakov, A. (1998). *Wavelets, Approximation, and Statistical Applications*. Springer.
- [30] d'Haultfoeuille, X. (2011). On the completeness condition for nonparametric instrumental problems. *Econometric Theory* **27** 460-471.
- [31] Hofinger, A. and Pikkarainen, H.K. (2007). Convergence rate for the Bayesian ap-

- proach to linear inverse problems. *Inverse Problems* **23** 2469-2484.
- [32] Hofinger, A. and Pikkarainen, H.K. (2009). Convergence rates for linear inverse problems in the presence of an additive normal noise. *Stoch. Anal. Appl.* **27** 240-257.
- [33] Horowitz, J.L. (2011). Applied nonparametric instrumental variables estimation. *Econometrica* **79** 347-394.
- [34] Horowitz, J.L. (2012). Specification testing in nonparametric instrumental variables estimation. *J. Econometrics* **167** 383-396.
- [35] Horowitz, J.L. and Lee, S. (2007). Nonparametric instrumental variables estimation of a quantile regression model. *Econometrica* **75** 1191-1208.
- [36] Hu, Y. and Shiu, J.-L. (2011). Nonparametric identification using instrumental variables: sufficient conditions for completeness. Preprint.
- [37] Jiang, W. and Tanner, M.A. (2008). Gibbs posterior for variable selection in high-dimensional classification and data mining. *Ann. Statist.* **36** 2207-2231.
- [38] Johnstone, I.M. (2011). *Gaussian Estimation: Sequence and Multiresolution Models*. Unpublished draft.
- [39] Kim, J. (2002). Limited information likelihood and Bayesian analysis. *J. Econometrics* **107** 175-193.
- [40] Kleijn, B.J.K. and van der Vaart, A.W. (2006). Misspecification in infinite-dimensional Bayesian statistics. *Ann. Statist.* **34** 837-877.
- [41] Knapik, B.T., Szabó, B.T., van der Vaart, A.W., and van Zanten, J.H. (2012). Bayes procedures for adaptive inference in inverse problems for the white noise model. arXiv:1209.3628.
- [42] Knapik, B.T., van der Vaart, A.W. and van Zanten, J.H. (2011). Bayesian inverse problems with Gaussian priors. *Ann. Statist.* **39** 2626-2657.
- [43] Knapik, B.T., van der Vaart, A.W. and van Zanten, J.H. (2013). Bayesian recovery of the initial condition for the heat equation. *Comm. Statist. Theory Methods* **42** 1294-1313.
- [44] Kovchegov, Y. and Yildiz, N. (2012). Identification via completeness for discrete covariates and orthogonal polynomials. Preprint.
- [45] Kress, R. (1999). *Linear Integral Equations*. Second Edition. Springer.
- [46] Liao, Y. and Jiang, W. (2011). Posterior consistency of nonparametric conditional moment restricted models. *Ann. Statist.* **39** 3003-3031.
- [47] Liu, J.S., Tian, L. and Wei, L.J. (2007). Implementation of estimating-function based inference procedures with Markov Chain Monte Carlo samplers. *J. Amer. Stat. Assoc.* **102** 881-888.
- [48] Mallat, S. (2009). *A Wavelet Tour of Signal Processing*. Third Edition. Academic Press.
- [49] Massart, P. (2000). About the constants in Talagrand's concentration inequalities for empirical processes. *Ann. Probab.* **28** 863-884.
- [50] Newey, W. and Powell, J. (2003). Instrumental variables estimation of non-parametric models. *Econometrica* **71** 1565-1578.
- [51] Rudelson, M. (1999). Random vectors in the isotropic position. *J. Functional Anal.* **164** 60-72.
- [52] Santos, A. (2012). Inference in nonparametric instrumental variables with partial identification. *Econometrica* **80** 213-275.
- [53] Schennach, S. (2005). Bayesian exponentially tilted empirical likelihood. *Biometrika* **92** 31-46.
- [54] Scricciolo, C. (2006). Convergence rates for Bayesian density estimation of infinite-dimensional exponential families. *Ann. Statist.* **34** 2897-2920.
- [55] Severini, T. and Tripathi, G. (2006). Some identification issues in nonparametric

- linear models with endogenous regressors. *Econometric Theory* **22** 258-278.
- [56] Shen, X. and Wasserman, L. (2001). Rates of convergence of posterior distributions. *Ann. Statist.* **29** 687-714.
- [57] Stuart, A.M. (2010). Inverse problems: a Bayesian perspective. *Acta Numerica* **19** 451-559.
- [58] Talagrand, M. (1996). New concentration inequalities in product spaces. *Invent. Math.* **126** 505-563.
- [59] van der Vaart, A.W. (1998). *Asymptotic Statistics*. Cambridge University Press.
- [60] van der Vaart, A.W. and Wellner, J.A. (1996). *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer.
- [61] Yosida, K. (1980). *Functional Analysis*. Springer.
- [62] Zhang, T. (2006a). From  $\epsilon$ -entropy to KL entropy: analysis of minimum information complexity density estimation. *Ann. Statist.* **34** 2180-2210.
- [63] Zhang, T. (2006b). Information theoretical upper and lower bounds for statistical estimation. *IEEE Trans. Inform. Theory* **54** 1307-1321.

## APPENDIX A: CDV WAVELET BASES AND BESOV SPACES

**A.1. Wavelet bases for  $L_2[0, 1]$ .** We review wavelet theory on the compact interval  $[0, 1]$ . We refer the reader to [29], [48] and Johnstone [38, Chapter 7 and Appendix B] as useful general references on wavelet theory in the statistical (and signal processing) context.

Let  $(\varphi, \psi)$  be a Daubechies pair of the scaling function and wavelet of a multiresolution analysis of the space  $L_2(\mathbb{R})$  of order  $N$ , with  $\psi$  having  $N$  vanishing moments and support contained in  $[-N + 1, N]$ , and  $\varphi$  having support contained in  $[0, 2N - 1]$  [see 29, Remark 7.1]. We translate  $\varphi$  so that its support is contained in  $[-N + 1, N]$ . Define

$$\varphi_{jk}(x) = 2^{j/2}\varphi(2^j x - k), \quad \psi_{jk}(x) = 2^{j/2}\psi(2^j x - k).$$

Then, for any fixed  $J_0 \geq 0$ , it is known that  $\{\varphi_{J_0 k}, \psi_{jk}, j \geq J_0, k \in \mathbb{Z}\}$  forms an orthonormal basis for  $L_2(\mathbb{R})$ . However, we need an orthonormal basis for  $L_2[0, 1]$ . From the Daubechies pair  $(\varphi, \psi)$ , we wish to construct an orthonormal basis for  $L_2[0, 1]$ . The construction here is based on Cohen et al. [15, Section 4]. See also Chapter 7.5 of [48] for wavelet bases on  $[0, 1]$ .

Take a fixed resolution level  $j$  such that  $2^j \geq 2N$ . For  $k = N, \dots, 2^j - N - 1$ ,  $\varphi_{jk}$  are supported in  $[0, 1]$  and left unchanged:  $\varphi_{jk}^{\text{int}}(x) = \varphi_{jk}(x)$  for  $x \in [0, 1]$ . At boundaries,  $k = 0, \dots, N - 1$ , construct suitable functions  $\varphi_k^L$  with support  $[0, N + k]$  and  $\varphi_k^R$  with support  $[-N - k, 0]$ , and define

$$\varphi_{jk}^{\text{int}}(x) = 2^{j/2}\varphi_k^L(2^j x), \quad \varphi_{j, 2^j - k - 1}^{\text{int}}(x) = 2^{j/2}\varphi_k^R(2^j(x - 1)), \quad x \in [0, 1].$$

Note that both  $\varphi_k^L$  and  $\varphi_k^R$  have the same smoothness as  $\varphi$ . Define the multiresolution spaces  $V_j = \text{span}\{\varphi_{jk}^{\text{int}}, k = 0, \dots, 2^j - 1\}$ , which satisfy the following properties (i)  $\dim(V_j) = 2^j$ ; (ii)  $V_j \subset V_{j+1}$ ; (iii) each  $V_j$  contains all polynomials of order at most  $N - 1$ .

Turning to the wavelet spaces, define  $W_j$  by the orthogonal complement of  $V_j$  in  $V_{j+1}$ . Starting from the Daubechies wavelet  $\psi$ , construct  $\psi_{jk}^{\text{int}}$  similarly to  $\varphi_{jk}^{\text{int}}$ . Then, we have  $W_j = \text{span}\{\psi_{jk}^{\text{int}}, k = 0, \dots, 2^j - 1\}$ , and for any  $J_0 \geq 1$  with  $2^{J_0} \geq 2N$  and  $J > J_0$ ,

$$V_J = V_{J_0} \bigoplus_{j \geq J_0}^{J-1} W_j, \quad L_2[0, 1] = V_{J_0} \bigoplus_{j \geq J_0} W_j.$$

Therefore,  $\{\varphi_{J_0 k}^{\text{int}}\}_{k=0}^{2^{J_0}-1} \cup \{\psi_{jk}^{\text{int}}, j \geq J_0, k = 0, \dots, 2^j - 1\}$  forms an orthonormal basis for  $L_2[0, 1]$  [see Section 4 of 15, for formal proofs of these results]

DEFINITION 1. Call the so-constructed basis  $\{\varphi_{J_0 k}^{\text{int}}\}_{k=0}^{2^{J_0}-1} \cup \{\psi_{jk}^{\text{int}}, j \geq J_0, k = 0, \dots, 2^j - 1\}$  the CDV (Cohen-Daubechies-Vial) wavelet basis for  $L_2[0, 1]$  generated from the Daubechies pair  $(\varphi, \psi)$ . If  $(\varphi, \psi)$  is  $S$ -regular, i.e., if  $(\varphi, \psi)$  are  $S$ -times continuously differentiable, then call the so-generated CDV wavelet basis  $S$ -regular.

REMARK 10. For any given positive integer  $S$ , there is an  $S$ -regular Daubechies pair  $(\varphi, \psi)$  by taking the order  $N$  sufficiently large [see 29, Remark 7.1].

**A.2. Besov spaces.** We recall the definition of Besov spaces.

DEFINITION 2. Let  $0 < s < S, s \in \mathbb{R}, S \in \mathbb{N}$  and  $1 \leq p, q \leq \infty$ . Let  $\{\varphi_{J_0 k}^{\text{int}}\}_{k=0}^{2^{J_0}-1} \cup \{\psi_{jk}^{\text{int}}, j \geq J_0, k = 0, \dots, 2^j - 1\}$  be an  $S$ -regular CDV wavelet basis for  $L_2[0, 1]$ . Let

$$\varphi_{J_0 k}^{\text{int}}(f) = \int_0^1 f(x) \varphi_{J_0 k}^{\text{int}}(x) dx, \quad \psi_{jk}^{\text{int}}(f) = \int_0^1 f(x) \psi_{jk}^{\text{int}}(x) dx.$$

Then the Besov space  $B_{p,q}^s$  is defined by the set of functions  $\{f \in L_2[0, 1] : \|f\|_{s,p,q} < \infty\}$ , where

$$\|f\|_{s,p,q} := \left( \sum_{0 \leq k \leq 2^{J_0}-1} |\varphi_{J_0 k}^{\text{int}}(f)|^p \right)^{1/p} + \left( \sum_{j \geq J_0} \left( 2^{j(s+1/2-1/p)} \left( \sum_{0 \leq k \leq 2^j-1} |\psi_{jk}^{\text{int}}(f)|^p \right)^{1/p} \right)^q \right)^{1/q},$$

with the obvious modification in case  $p = \infty$  or  $q = \infty$ .

REMARK 11. Besov spaces cover commonly used smooth function spaces. For example,  $B_{\infty,\infty}^s$  is equal to the Hölder-Zygmund space, which coincides with the classical Hölder space for non-integer  $s$ . For integer  $s$ , they do not coincide but the Hölder-Zygmund space contains the classical Hölder space. Moreover,  $B_{2,2}^s$  is equal to the classical  $L_2$ -Sobolev space when  $s$  is an integer. See [38], Appendix B.

## APPENDIX B: PROOFS OF LEMMAS 1-3

PROOF OF LEMMA 1. For part (ii), the lower bound on  $s_{\min}(\mathbb{E}[\phi^J(W)^{\otimes 2}])$  follows from Assumption 1 (iii); the upper bounds on  $s_{\max}(\mathbb{E}[\phi^J(W)^{\otimes 2}])$  and

$s_{\max}(\mathbb{E}[\phi^J(W)\phi^J(X)^T])$  follow from Assumption 1 (i) and the fact that  $\{\phi_l\}$  is an orthonormal basis of  $L_2[0, 1]$  (see (8)). Part (iii) follows from Rudelson's [51] inequality and (i). For the reader's convenience, we state Rudelson's inequality in Appendix E. For Part (v), we first note that, by (iii) and Weyl's perturbation theorem [5, Problem III.6.13],  $s_{\min}(\hat{\Phi}_{WX}) \geq \tau_{J_n} - O_P(\sqrt{J_n 2^{J_n}/n})$ . Since now  $\sqrt{J_n 2^{J_n}/n} = o(\tau_{J_n})$ , we have  $s_{\min}(\hat{\Phi}_{WX}) \geq (1 - o_P(1))\tau_{J_n}$ . For the proof of (i), denote by  $N$  the order of the Daubechies pair  $(\varphi, \psi)$  generating the CDV wavelet basis  $\{\phi_l, l \geq 1\}$ . Then, for each  $x \in [0, 1]$  and each  $j \geq J_0$ , the number of nonzero elements in  $\phi_{2^j+1}(x), \dots, \phi_{2^{j+1}}(x)$  is bounded by some constant depending only on  $N$ , and each  $\phi_{2^j+k}(x)$  is bounded by some constant (depending only on  $\psi$ ) times  $2^{j/2}$  for all  $k = 1, \dots, 2^j$ . Similarly,  $\phi_1, \dots, \phi_{2^{J_0}}$  are uniformly bounded. Therefore, there exists a constant  $D$  depending only on  $(\varphi, \psi)$  such that  $\|\phi^J(x)\|_{\ell^2}^2 \leq D(2^{J_0} + \sum_{j=J_0}^{J-1} 2^j) = D2^J$  for all  $x \in [0, 1]$ .

Finally, we wish to show Part (iv). First, observe that  $\|\mathbb{E}_n[\phi^{J_n}(W_i)R_i]\|_{\ell^2}^2 \leq 2\|\mathbb{E}_n[\phi^{J_n}(W_i)R_i] - \mathbb{E}[\phi^{J_n}(W)R]\|_{\ell^2}^2 + 2\|\mathbb{E}[\phi^{J_n}(W)R]\|_{\ell^2}^2$ . By a simple moment calculation, the first term is  $O_P(2^{J_n}/n)$ . For the second term, by Assumptions 3 and 4 (ii),

$$\begin{aligned} \|\mathbb{E}[\phi^{J_n}(W)R]\|_{\ell^2}^2 &= \|\mathbb{E}[\phi^{J_n}(W)(g_0 - P_J g_0)(X)]\|_{\ell^2}^2 \\ &\lesssim \tau_{J_n}^2 \|g_0 - P_{J_n} g_0\|^2 \\ &\lesssim \tau_{J_n}^2 2^{-2J_n s}. \end{aligned}$$

This completes the proof.  $\square$

PROOF OF LEMMA 2. Step 1. We first show that  $|\Sigma_n| = 1 + o(1)$  ( $|\Sigma_n|$  denotes the determinant of  $\Sigma_n$ ). Let  $\lambda_{\min, n}$  and  $\lambda_{\max, n}$  denote the minimum and maximum eigenvalues of  $\Sigma_n$ , respectively. Then, by Weyl's perturbation theorem,  $1 - o(k_n^{-1}) \leq \lambda_{\min, n} \leq \lambda_{\max, n} \leq 1 + o(k_n^{-1})$ , so that  $(1 - o(k_n^{-1}))^{k_n} = \lambda_{\min, n}^{k_n} \leq |\Sigma_n| \leq \lambda_{\max, n}^{k_n} = (1 + o(k_n^{-1}))^{k_n}$ . Here both sides converge to 1.

Step 2. By Step 1, we have

$$\begin{aligned} &\int |dN(0, \Sigma_n)(x) - dN(0, I_{k_n})(x)| dx \\ &= \frac{1}{(2\pi)^{k_n/2}} \int \left| \frac{1}{|\Sigma_n|^{1/2}} e^{-x^T \Sigma_n^{-1} x/2} - e^{-x^T x/2} \right| dx \\ &\leq \left| \frac{1}{|\Sigma_n|^{1/2}} - 1 \right| + \frac{1}{(2\pi)^{k_n/2} |\Sigma_n|^{1/2}} \int |e^{-x^T \Sigma_n^{-1} x/2} - e^{-x^T x/2}| dx \\ &\leq o(1) + \frac{1}{(2\pi)^{k_n/2} (1 + o(1))} \int e^{-x^T x/2} |e^{-x^T (\Sigma_n^{-1} - I_{k_n}) x/2} - 1| dx \end{aligned}$$

By assumption, we have  $\epsilon_n := \|\Sigma_n^{-1} - I_{k_n}\|_{\text{op}} \leq \|\Sigma_n^{-1}\|_{\text{op}} \|I_{k_n} - \Sigma_n\|_{\text{op}} = o(k_n^{-1})$ . Now,  $|e^{-x^T(\Sigma_n^{-1} - I_{k_n})x/2} - 1| \leq e^{\epsilon_n x^T x/2} - e^{-\epsilon_n x^T x/2}$ . By a direct calculation, the conclusion follows from the fact that  $(1 \pm \epsilon_n)^{k_n} = 1 + o(1)$ .  $\square$

PROOF OF LEMMA 3. The first assertion follows from the assumption. Suppose now that  $\hat{A}_n$  is non-singular. Then,  $\hat{A}_n^{-1}A_n = (\hat{A}_n - A_n + A_n)^{-1}A_n = (A_n^{-1}\hat{A}_n - I_{k_n} + I_{k_n})^{-1}$ . Here,  $A_n^{-1}\hat{A}_n - I_{k_n} = A_n^{-1}(\hat{A}_n - A_n)$ , so that  $\|A_n^{-1}\hat{A}_n - I_{k_n}\|_{\text{op}} \leq \|A_n^{-1}\|_{\text{op}} \|\hat{A}_n - A_n\|_{\text{op}} = s_{\min}^{-1}(A_n) \|\hat{A}_n - A_n\|_{\text{op}} = O_P(\epsilon_n^{-1}\delta_n)$ . Let  $\hat{\Delta} = I_{k_n} - A_n^{-1}\hat{A}_n$ . Then,  $\hat{A}_n^{-1}A_n = (I_{k_n} - \hat{\Delta})^{-1} = I_{k_n} + \sum_{m=1}^{\infty} \hat{\Delta}^m$  (Neumann series). Therefore, we conclude that  $\|\hat{A}_n^{-1}A_n - I_{k_n}\|_{\text{op}} = \|\sum_{m=1}^{\infty} \hat{\Delta}^m\|_{\text{op}} \leq \sum_{m=1}^{\infty} \|\hat{\Delta}\|_{\text{op}}^m = \|\hat{\Delta}\|_{\text{op}} \cdot \sum_{m=0}^{\infty} \|\hat{\Delta}\|_{\text{op}}^m = O_P(\epsilon_n^{-1}\delta_n)$ .  $\square$

### APPENDIX C: PROOF OF THEOREM 3

For the notational convenience, define

$$\mathbb{E}_{\Pi_n}[\cdot | \mathcal{D}_n] := \int \cdot \Pi_n(dg | \mathcal{D}_n), \quad \mathbb{E}_{\tilde{\Pi}_n}[\cdot | \mathcal{D}_n] := \int \cdot \tilde{\Pi}_n(db^{J_n} | \mathcal{D}_n).$$

PROOF OF THEOREM 3. Define the event

$$\mathcal{E}_{3n} = \{\mathcal{D}_n : \hat{\Phi}_{WX} \text{ and } \hat{\Phi}_{WW} \text{ are non-singular}\}.$$

Then, by Lemma 1,  $\mathbb{P}\{1(\mathcal{E}_{3n}) = 1\} = \mathbb{P}(\mathcal{E}_{3n}) \rightarrow 1$ . Suppose that  $1(\mathcal{E}_{3n}) = 1$ . Then, by (16),  $\ell_{b^{J_n}}(\mathcal{D}_n)$  defined in the proof of Proposition 4 is bounded from below by

$$\hat{c} \|b^{J_n}\|_{\ell^2}^2 + \text{a term independent of } b^{J_n},$$

for some positive random variable  $\hat{c}$ . Hence, the integral  $\mathbb{E}_{\tilde{\Pi}_n}[\|b^{J_n}\|_{\ell^2} | \mathcal{D}_n]$  is finite as soon as  $1(\mathcal{E}_{3n}) = 1$ . This proves the first assertion.

In what follows, we wish to prove the convergence rate result (14). First of all, by the triangle inequality and Jensen's inequality,

$$\begin{aligned} 1(\mathcal{E}_{3n}) \|\hat{g}_{QB} - g_0\| &\leq 1(\mathcal{E}_{3n}) \|\hat{g}_{QB} - P_{J_n}g_0\| + \|g_0 - P_{J_n}g_0\| \\ &= 1(\mathcal{E}_{3n}) \|\mathbb{E}_{\Pi_n}[g - P_{J_n}g_0 | \mathcal{D}_n]\| + \|g_0 - P_{J_n}g_0\| \\ &= 1(\mathcal{E}_{3n}) \|\mathbb{E}_{\tilde{\Pi}_n}[b^{J_n} - b_0^{J_n} | \mathcal{D}_n]\|_{\ell^2} + \|g_0 - P_{J_n}g_0\| \\ &\leq 1(\mathcal{E}_{3n}) \mathbb{E}_{\tilde{\Pi}_n}[\|b^{J_n} - b_0^{J_n}\|_{\ell^2} | \mathcal{D}_n] + \|g_0 - P_{J_n}g_0\|. \end{aligned}$$

Since  $\|g_0 - P_{J_n}g_0\| = O(2^{-J_n s})$ , it suffices to show that there exists a constant  $D > 0$  such that for every  $M_n \rightarrow \infty$ ,

$$\begin{aligned} \mathbb{P}\left[1(\mathcal{E}_{3n}) \mathbb{E}_{\tilde{\Pi}_n}[\|b^{J_n} - b_0^{J_n}\|_{\ell^2} | \mathcal{D}_n] \right. \\ \left. \leq D \max\{2^{-J_n s}, \tau_{J_n}^{-1} \sqrt{2^{J_n}/n}, \tau_{J_n}^{-1} \epsilon_n \varrho_n M_n\} \right] \rightarrow 1. \end{aligned}$$

Let  $\pi_n^*(\theta^{J_n} | \mathcal{D}_n)$  be the (random) density defined in the proof of Theorem 1. Note that  $\pi_n^*(\theta^{J_n} | \mathcal{D}_n)$  is well-defined as soon as  $1(\mathcal{E}_{3n}) = 1$ . Let  $\delta_n := \epsilon_n + \tau_{J_n} 2^{-J_n s}$ . Then we have:

LEMMA 5. *There exists a constant  $c_1 > 0$  such that for every sequence  $M_n \rightarrow \infty$  with  $M_n = o(L_n)$ ,*

$$\mathbb{P} \left\{ 1(\mathcal{E}_{3n}) \int \|\theta^{J_n}\|_{\ell^2} \cdot |\pi_n^*(\theta^{J_n} | \mathcal{D}_n) - dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})| d\theta^{J_n} \leq e^{-c_1 M_n n \delta_n^2} + M_n \sqrt{n} \delta_n \varrho_n \right\} \rightarrow 1,$$

where  $\Delta_n := \sqrt{n} \mathbb{E}_n[\phi^{J_n}(W_i) R_i]$ .

We defer the proof of Lemma 5 after the proof of this theorem. Note that

$$\begin{aligned} 1(\mathcal{E}_{3n}) & \left[ \int \|\theta^{J_n}\|_{\ell^2} dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n}) d\theta^{J_n} \right]^2 \\ & \leq 1(\mathcal{E}_{3n}) \int \|\theta^{J_n}\|_{\ell^2}^2 dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n}) d\theta^{J_n} \\ & \leq \|\Delta_n\|_{\ell^2}^2 + \text{tr}(\hat{\Phi}_{WW}). \end{aligned}$$

By the proof of Theorem 2, there exists a constant  $D_1 > 0$  such that  $\mathbb{P}\{\|\Delta_n\|_{\ell^2}^2 + \text{tr}(\hat{\Phi}_{WW}) \leq D_1(n\tau_{J_n}^2 2^{-2J_n s} + 2^{J_n})\} \rightarrow 1$ . Hence for every sequence  $M_n \rightarrow \infty$  with  $M_n = o(L_n)$ , with probability approaching one,

$$\begin{aligned} & \sqrt{D_1(n\tau_{J_n}^2 2^{-2J_n s} + 2^{J_n})} + e^{-c_1 M_n n \delta_n^2} + M_n \sqrt{n} \delta_n \varrho_n \\ & \geq 1(\mathcal{E}_{3n}) \int \|\theta^{J_n}\|_{\ell^2} \pi_n^*(\theta^{J_n} | \mathcal{D}_n) \\ & = 1(\mathcal{E}_{3n}) \sqrt{n} \int \|\hat{\Phi}_{WX}(b^{J_n} - b_0^{J_n})\|_{\ell^2} \tilde{\pi}_n(b^{J_n} | \mathcal{D}_n) db^{J_n} \\ & \geq 1(\mathcal{E}_{3n}) \sqrt{n} s_{\min}(\hat{\Phi}_{WX}) \mathbb{E}_{\tilde{\Pi}_n}[\|b^{J_n} - b_0^{J_n}\|_{\ell^2} | \mathcal{D}_n]. \end{aligned}$$

Take  $M_n \rightarrow \infty$  sufficiently slowly such that  $\varrho_n M_n \rightarrow 0$ . Since the left side is then  $\lesssim \max\{\sqrt{n} \tau_{J_n} 2^{-J_n s}, \sqrt{2^{J_n}}, \sqrt{n} \epsilon_n \varrho_n M_n\}$ , there exists a constant  $D_2 > 0$  such that

$$\begin{aligned} & \mathbb{P} \left[ 1(\mathcal{E}_{3n}) s_{\min}(\hat{\Phi}_{WX}) \mathbb{E}_{\tilde{\Pi}_n}[\|b^{J_n} - b_0^{J_n}\|_{\ell^2} | \mathcal{D}_n] \leq D_2 \max\{\tau_{J_n} 2^{-J_n s}, \sqrt{2^{J_n}/n}, \epsilon_n \varrho_n M_n\} \right] \rightarrow 1. \end{aligned}$$

Finally, by Lemma 1,  $\mathbb{P}(s_{\min}(\hat{\Phi}_{WX}) \geq 0.5\tau_{J_n}) \rightarrow 1$ , by which we have

$$\begin{aligned} & \mathbb{P}\left[1(\mathcal{E}_{3n})\mathbb{E}_{\hat{\Pi}_n}[\|b^{J_n} - b_0^{J_n}\|_{\ell^2} \mid \mathcal{D}_n] \right. \\ & \quad \left. \leq 2D_2 \max\{2^{-J_n s}, \tau_{J_n}^{-1}\sqrt{2^{J_n}/n}, \tau_{J_n}^{-1}\epsilon_n \varrho_n M_n\}\right] \rightarrow 1. \end{aligned}$$

This leads to the desired conclusion (it is not difficult to see that the final expression holds for every sequence  $M_n \rightarrow \infty$ ).  $\square$

PROOF OF LEMMA 5. As before, we say that a sequence of random variables  $A_n$  is *eventually* bounded by another sequence of random variables  $B_n$  if  $\mathbb{P}(A_n \leq B_n) \rightarrow 1$ .

Take any  $M_n \rightarrow \infty$  with  $M_n = o(L_n)$ . Then,

$$\begin{aligned} & 1(\mathcal{E}_{3n}) \int \|\theta^{J_n}\|_{\ell^2} \cdot |\pi_n^*(\theta^{J_n} \mid \mathcal{D}_n) - dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})| d\theta^{J_n} \\ & \leq 1(\mathcal{E}_{3n}) \int_{\|\theta^{J_n}\|_{\ell^2} \leq M_n \sqrt{n}\delta_n} \|\theta^{J_n}\|_{\ell^2} \cdot |\pi_n^*(\theta^{J_n} \mid \mathcal{D}_n) - dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})| d\theta^{J_n} \\ & \quad + 1(\mathcal{E}_{3n}) \int_{\|\theta^{J_n}\|_{\ell^2} > M_n \sqrt{n}\delta_n} \|\theta^{J_n}\|_{\ell^2} \pi_n^*(\theta^{J_n} \mid \mathcal{D}_n) d\theta^{J_n} \\ & \quad + 1(\mathcal{E}_{3n}) \int_{\|\theta^{J_n}\|_{\ell^2} > M_n \sqrt{n}\delta_n} \|\theta^{J_n}\|_{\ell^2} dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n}) d\theta^{J_n} \\ & =: I + II + III. \end{aligned}$$

We divide the rest of the proof into three steps.

Step 1. Claim:  $\exists c_2 > 0$  such that  $\mathbb{P}(II \leq e^{-c_2 M_n^2 n \delta_n^2}) \rightarrow 1$ .

(Proof of Step 1): The assertion of Step 1 follows from the same line as in the proof of Proposition 4 by noting that for any  $c > 0$ ,  $x e^{-cx^2} \leq e^{-cx^2/2}$  for all  $x > 0$  sufficiently large. Hence the proof is omitted.

Step 2. Claim:  $\exists c_3 > 0$  such that  $\mathbb{P}(III \leq e^{-c_3 M_n n \delta_n^2}) \rightarrow 1$ .

(Proof of Step 2): By the Cauchy-Schwarz inequality, the square of  $III$  is bounded by  $\int \|\theta^{J_n}\|_{\ell^2}^2 dN(\Delta_n, \hat{\Phi}_{WW}) d\theta^{J_n} \int_{\|\theta^{J_n}\|_{\ell^2} > M_n \sqrt{n}\delta_n} dN(\Delta_n, \hat{\Phi}_{WW}) d\theta^{J_n}$ . Here the first integral is bounded by

$$\|\Delta_n\|_{\ell^2}^2 + \text{tr}(\hat{\Phi}_{WW}),$$

which is eventually bounded by  $D(n\tau_{J_n}^2 2^{-2J_n s} + 2^{J_n})$  for some constant  $D > 0$  by the proof of Theorem 2. On the other hand, by the proof of Theorem 1, the second integral is eventually bounded by  $\int_{\|\theta^{J_n}\|_{\ell^2} > \sqrt{M_n n \delta_n}} dN(0, I_{2^{J_n}}) d\theta^{J_n}$ . By Borell's inequality for Gaussian measures (see (23)), the last integral

is bounded by  $e^{-c'M_n n \delta_n^2}$  for some small constant  $c' > 0$ . Taking these together, we obtain the conclusion of Step 2 by choosing the constant  $c_3 > 0$  sufficiently small.

Step 3. Claim:  $\exists c_4 > 0$  such that  $\mathbb{P}(I \leq e^{-c_4 M_n^2 n \delta_n^2} + M_n \sqrt{n} \delta_n \varrho_n) \rightarrow 1$ .

(Proof of Step 3): Let  $\mathcal{C}_n := \{\theta^{J_n} \in \mathbb{R}^{2^{J_n}} : \|\theta^{J_n}\|_{\ell^2} \leq M_n \sqrt{n} \delta_n\}$ . Let  $\pi_{n, \mathcal{C}_n}^*(\theta^{J_n} | \mathcal{D}_n)$  and  $dN^{\mathcal{C}_n}(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})$  denote the probability densities obtained by first restricting  $\pi_n^*(\theta^{J_n} | \mathcal{D}_n)$  and  $dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})$  to the ball  $\mathcal{C}_n$  and then renormalizing, respectively. Then, abbreviating  $\pi_n^*(d\theta^{J_n} | \mathcal{D}_n)$  by  $\pi_n^*$ ,  $\pi_{n, \mathcal{C}_n}^*(d\theta^{J_n} | \mathcal{D}_n)$  by  $\pi_{n, \mathcal{C}_n}^*$ ,  $dN(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})$  by  $dN$ , and  $dN^{\mathcal{C}_n}(\Delta_n, \hat{\Phi}_{WW})(\theta^{J_n})$  by  $dN^{\mathcal{C}_n}$  we have

$$\begin{aligned} I &\leq 1(\mathcal{E}_{3n}) \int \|\theta^{J_n}\|_{\ell^2} \cdot |\pi_{n, \mathcal{C}_n}^* - dN^{\mathcal{C}_n}| \\ &\quad + 1(\mathcal{E}_{3n}) \int_{\|\theta^{J_n}\|_{\ell^2} \leq M_n \sqrt{n} \delta_n} \|\theta^{J_n}\|_{\ell^2} \cdot |\pi_{n, \mathcal{C}_n}^* - \pi_n^*| \\ &\quad + 1(\mathcal{E}_{3n}) \int_{\|\theta^{J_n}\|_{\ell^2} \leq M_n \sqrt{n} \delta_n} \|\theta^{J_n}\|_{\ell^2} \cdot |dN^{\mathcal{C}_n} - dN| \\ &=: IV + V + IV. \end{aligned}$$

By the proof of Theorem 1, the term  $IV$  is eventually bounded by

$$1(\mathcal{E}_{3n}) M_n \sqrt{n} \delta_n \int |\pi_{n, \mathcal{C}_n}^* - dN^{\mathcal{C}_n}| \leq M_n \sqrt{n} \delta_n \varrho_n.$$

For the term  $V$ , we have

$$\begin{aligned} V &\leq 1(\mathcal{E}_{3n}) M_n \sqrt{n} \delta_n \int_{\|\theta^{J_n}\|_{\ell^2} \leq M_n \sqrt{n} \delta_n} |\pi_{n, \mathcal{C}_n}^* - \pi_n^*| \\ &= 1(\mathcal{E}_{3n}) M_n \sqrt{n} \delta_n \int_{\|\theta^{J_n}\|_{\ell^2} > M_n \sqrt{n} \delta_n} \pi_n^*. \end{aligned}$$

By the proof of Proposition 4, there exists a constant  $c_5 > 0$  such that the integral on the right side is eventually bounded by  $e^{-c_5 M_n^2 n \delta_n^2}$ , so that  $\mathbb{P}(V \leq e^{-c_5 M_n^2 n \delta_n^2 / 2}) \rightarrow 1$ . Likewise, by Borell's inequality for Gaussian measures, there exists a constant  $c_6 > 0$  such that  $\mathbb{P}(VI \leq e^{-c_6 M_n n \delta_n^2}) \rightarrow 1$ . Taking these together, we obtain the conclusion of Step 3 by choosing the constant  $c_4 > 0$  sufficiently small.

Finally, Steps 1-3 lead to the conclusion of Lemma 5.  $\square$

## APPENDIX D: PROOFS FOR SECTION 4

PROOF OF PROPOSITION 2. We only consider the mildly ill-posed case. The proof for the severely ill-posed case is similar. For either case of product or isotropic priors, it suffices to check conditions P1) and P2) in Theorem 1. We shall do this with the choice  $\epsilon_n = \sqrt{2^{J_n}(\log n)/n} \sim (\log n)^{1/2} n^{-(r+s)/(2r+2s+1)}$ .

Case of product priors: Let  $c_{\min} := \min_{x \in [-A, A]} q(x) > 0$ . Since  $\|b^{J_n} - b_0^{J_n}\|_{\ell^2}^2 = \sum_{l=1}^{2^{J_n}} (b_l - b_{0l})^2 \leq 2^{J_n} \max_{1 \leq l \leq 2^{J_n}} (b_l - b_{0l})^2$ , we have

$$\begin{aligned} \tilde{\Pi}_n(b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\ell^2} \leq \epsilon_n) &\geq \tilde{\Pi}_n \left( b^{J_n} : \max_{1 \leq l \leq 2^{J_n}} |b_l - b_{0l}| \leq \epsilon_n / \sqrt{2^{J_n}} \right) \\ &\geq \prod_{l=1}^{2^{J_n}} \tilde{\Pi}_n(b^{J_n} : |b_l - b_{0l}| \leq \epsilon_n / \sqrt{2^{J_n}}). \end{aligned}$$

Since  $\exists \epsilon \in (0, A)$ ,  $b_{0l} \in [-A + \epsilon, A - \epsilon]$  for all  $l \geq 1$ , for all  $n$  sufficiently large, the last expression is bounded from below by

$$\left( \frac{c_{\min} \epsilon_n}{\sqrt{2^{J_n}}} \right)^{2^{J_n}} = e^{-2^{J_n} \log(\sqrt{2^{J_n}} / (c_{\min} \epsilon_n))} \geq e^{-C n \epsilon_n^2},$$

where  $C > 0$  is a sufficiently large constant, which verifies condition P1).

Second, with this  $\epsilon_n$ ,  $\gamma_n$  in condition P2) is  $\sim (\log n)^{1/2} n^{-s/(2r+2s+1)}$ . Let, say,  $L_n \sim (\log n)^{1/2}$  so that  $L_n \gamma_n \sim (\log n) n^{-s/(2r+2s+1)}$ . Then,  $\{b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\ell^2} \leq L_n \gamma_n\} \subset [-A, A]^{2^{J_n}}$  for all  $n$  sufficiently large, so that  $\tilde{\pi}_n(b^{J_n}) = \prod_{l=1}^{2^{J_n}} q(b_l)$  is positive for all  $\|b^{J_n} - b_0^{J_n}\|_{\ell^2} \leq L_n \gamma_n$ . Let  $\|b^{J_n}\|_{\ell^2} \leq L_n \gamma_n$  and  $\|b^{J_n}\|_{\ell^2} \leq L_n \gamma_n$ . Then,

$$\begin{aligned} \frac{\tilde{\pi}_n(b_0^{J_n} + b^{J_n})}{\tilde{\pi}_n(b_0^{J_n} + \tilde{b}^{J_n})} &= \exp \left[ \sum_{l=1}^{2^{J_n}} \{ \log q(b_{0l} + b_l) - \log q(b_{0l} + \tilde{b}_l) \} \right] \\ &\leq \exp \left\{ L \sum_{l=1}^{2^{J_n}} |b_l - \tilde{b}_l| \right\} \leq \exp \left\{ L \sqrt{2^{J_n}} \|b^{J_n} - \tilde{b}^{J_n}\|_{\ell^2} \right\} \\ &\leq e^{2L \sqrt{2^{J_n}} L_n \gamma_n} = e^{o(1)}, \end{aligned}$$

where the last step is due to  $s > 1/2$ . Likewise, we have

$$\frac{\tilde{\pi}_n(b_0^{J_n} + b^{J_n})}{\tilde{\pi}_n(b_0^{J_n} + \tilde{b}^{J_n})} \geq e^{-2L \sqrt{2^{J_n}} L_n \gamma_n} = e^{-o(1)}.$$

Therefore, condition P2) is verified.

Case of isotropic priors: Let  $c_{\min} := \min_{x \in [0, A]} r(x) > 0$ . Then, for all  $n$  sufficiently large,

$$\begin{aligned} \tilde{\Pi}_n(b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\ell^2} \leq \epsilon_n) &= \frac{\int_{\|b^{J_n} - b_0^{J_n}\|_{\ell^2} \leq \epsilon_n} r(\|b^{J_n}\|_{\ell^2}) db^{J_n}}{\int r(\|b^{J_n}\|_{\ell^2}) db^{J_n}} \\ &= \frac{\int_{\|b^{J_n}\|_{\ell^2} \leq \epsilon_n} r(\|b^{J_n} + b_0^{J_n}\|_{\ell^2}) db^{J_n}}{\int r(\|b^{J_n}\|_{\ell^2}) db^{J_n}} \geq \frac{c_{\min} \int_{\|b^{J_n}\|_{\ell^2} \leq \epsilon_n} db^{J_n}}{\int r(\|b^{J_n}\|_{\ell^2}) db^{J_n}} \\ &= c_{\min} \frac{\int_{x \in [0, \epsilon_n]} x^{2^{J_n} - 1} dx}{\int_0^\infty x^{2^{J_n} - 1} r(x) dx} \geq c_{\min} \left( \frac{\epsilon_n}{2^{J_n}} \right)^{2^{J_n}} \times e^{-c'' 2^{J_n} \log(2^{J_n})} \\ &= c_{\min} e^{-2^{J_n} \log(2^{J_n} / \epsilon_n) - c'' 2^{J_n} \log(2^{J_n})} \geq e^{-C n \epsilon_n^2}, \end{aligned}$$

where  $C > 0$  is a sufficiently large constant, which verifies condition P1).

Second, with this  $\epsilon_n$ ,  $\gamma_n$  in condition P2) is  $\sim (\log n)^{1/2} n^{-s/(2r+2s+1)}$ . Let  $L_n \sim (\log n)^{1/2}$  so that  $L_n \gamma_n \sim (\log n) n^{-s/(2r+2s+1)}$ . Since  $\|b_0^{J_n}\|_{\ell^2} \leq \|g_0\| < A$  and  $L_n \gamma_n \rightarrow 0$ ,  $\{b^{J_n} : \|b^{J_n} - b_0^{J_n}\|_{\ell^2} \leq L_n \gamma_n\} \subset \{b^{J_n} : \|b^{J_n}\|_{\ell^2} \leq A\}$  for all  $n$  sufficiently large, so that  $\tilde{\pi}_n(b^{J_n}) \propto r(\|b^{J_n}\|_{\ell^2})$  is positive for all  $\|b^{J_n} - b_0^{J_n}\|_{\ell^2} \leq L_n \gamma_n$ . Let  $\|b^{J_n}\|_{\ell^2} \leq L_n \gamma_n$  and  $\|\tilde{b}^{J_n}\|_{\ell^2} \leq L_n \gamma_n$ . Then by Parseval's identity,

$$\|b_0^{J_n} + b^{J_n}\|_{\ell^2} \leq \|b_0^{J_n}\|_{\ell^2} + L_n \gamma_n \rightarrow \|g_0\|,$$

and likewise we have

$$\|b_0^{J_n} + b^{J_n}\|_{\ell^2} \geq \|b_0^{J_n}\|_{\ell^2} - L_n \gamma_n \rightarrow \|g_0\|.$$

Therefore, we conclude that, uniformly in  $\|b^{J_n}\|_{\ell^2} \leq L_n \gamma_n$ ,  $\|\tilde{b}^{J_n}\|_{\ell^2} \leq L_n \gamma_n$ ,

$$\frac{\tilde{\pi}_n(b_0^{J_n} + b^{J_n})}{\tilde{\pi}_n(b_0^{J_n} + \tilde{b}^{J_n})} = \frac{r(\|b_0^{J_n} + b^{J_n}\|_{\ell^2})}{r(\|b_0^{J_n} + \tilde{b}^{J_n}\|_{\ell^2})} \rightarrow \frac{r(\|g_0\|)}{r(\|g_0\|)} = 1.$$

Hence condition P2) is verified.  $\square$

PROOF OF PROPOSITION 3. Given the proof of Proposition 2 and the discussion following Theorem 3, it is sufficient to verify that  $\varrho_n$  is  $o((\log n)^{-1/2})$ . However, this is readily verified by tracking the proof of Proposition 2.  $\square$

## APPENDIX E: TECHNICAL TOOLS

We state here Rudelson's inequality for the reader's convenience.

THEOREM 4 (Rudelson's [51] inequality). *Let  $Z_1, \dots, Z_n$  be i.i.d. random vectors in  $\mathbb{R}^k$  with  $\Sigma := E[Z_1^{\otimes 2}]$ . Then for every  $k \geq e^2$ ,*

$$\mathbb{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n Z_i^{\otimes 2} - \Sigma \right\|_{\text{op}} \right] \leq \max\{\|\Sigma\|_{\text{op}}^{1/2} \delta, \delta^2\}, \quad \delta = D \sqrt{\frac{\log k}{n} \mathbb{E}[\max_{1 \leq i \leq n} \|Z_i\|_{\ell^2}^2]},$$

where  $D$  is a universal constant.

Rudelson's inequality implies the following corollary useful in our application.

COROLLARY 1. *Let  $(X_1, Y_1^T)^T, \dots, (X_n, Y_n^T)^T$  be i.i.d. random vectors with  $X_i \in \mathbb{R}^{k_1}, Y_i \in \mathbb{R}^{k_2}$ , and  $k_1 + k_2 \geq e^2$ . Let  $\Sigma_X := \mathbb{E}[X_1^{\otimes 2}], \Sigma_Y := \mathbb{E}[Y_1^{\otimes 2}]$  and  $\Sigma_{XY} := \mathbb{E}[X_1 Y_1^T]$ . Suppose that there exists a finite number  $m$  such that  $\mathbb{E}[\max_{1 \leq i \leq n} \|X_i\|_{\ell^2}^2] \vee \mathbb{E}[\max_{1 \leq i \leq n} \|Y_i\|_{\ell^2}^2] \leq m$ . Then*

$$\mathbb{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n X_i Y_i^T - \Sigma_{XY} \right\|_{\text{op}} \right] \leq \max\{(\|\Sigma_X\|_{\text{op}}^{1/2} \vee \|\Sigma_Y\|_{\text{op}}^{1/2}) \delta, \delta^2\},$$

with  $\delta = D \sqrt{\frac{m \log(k_1 \vee k_2)}{n}}$ ,

where  $D$  is a universal constant.

PROOF. Let  $Z_i = (X_i, Y_i^T)^T$ , and apply Rudelson's inequality to  $Z_1, \dots, Z_n$ . Note that by the variational characterization of the operator norm, we have  $\|n^{-1} \sum_{i=1}^n X_i Y_i^T - \Sigma_{XY}\|_{\text{op}} \leq \|n^{-1} \sum_{i=1}^n Z_i^{\otimes 2} - \mathbb{E}[Z_1^{\otimes 2}]\|_{\text{op}}$ , and by the Cauchy-Schwarz inequality,  $\|\mathbb{E}[Z_1^{\otimes 2}]\|_{\text{op}} \leq 2\|\Sigma_X\|_{\text{op}} + 2\|\Sigma_Y\|_{\text{op}}$ .  $\square$

Lastly, we shall recall Talagrand's [58] concentration inequality for general empirical processes. The following version is due to [49]. Here, for a generic class  $\mathcal{F}$  of measurable functions on some measurable space  $\mathcal{X}$ , we say that  $\mathcal{F}$  is pointwise measurable if there exists a countable subclass  $\mathcal{G} \subset \mathcal{F}$  such that for any  $f \in \mathcal{F}$ , there exists a sequence  $\{g_m\} \subset \mathcal{G}$  with  $g_m(x) \rightarrow f(x)$  for every  $x \in \mathcal{X}$ . See Chapter 2.3 of [60].

THEOREM 5 (Massart's form of Talagrand's inequality). *Let  $\xi_1, \dots, \xi_n$  be i.i.d. random variables taking values in some measurable space  $\mathcal{X}$ . Let  $\mathcal{F}$  be a pointwise measurable class of functions on  $\mathcal{X}$  such that  $\mathbb{E}[f(\xi_1)] = 0$  for all  $f \in \mathcal{F}$  and  $\sup_{f \in \mathcal{F}} \sup_{x \in \mathcal{X}} |f(x)| \leq B$  for some constant  $B > 0$ .*

Let  $\sigma^2$  be any positive constant such that  $\sigma^2 \geq \sup_{f \in \mathcal{F}} \mathbb{E}[f^2(\xi_1)]$ . Let  $Z := \sup_{f \in \mathcal{F}} |\sum_{i=1}^n f(\xi_i)|$ . Then for every  $x > 0$ ,

$$\mathbb{P}\{Z \geq C(\mathbb{E}[Z] + \sigma\sqrt{nx} + Bx)\} \leq e^{-x},$$

where  $C > 0$  is a universal constant.

GRADUATE SCHOOL OF ECONOMICS, UNIVERSITY OF TOKYO  
7-3-1 HONGO, BUNKYO-KU, TOKYO 113-0033, JAPAN.  
E-MAIL: [kkato@e.u-tokyo.ac.jp](mailto:kkato@e.u-tokyo.ac.jp)