

DATA-DRIVEN TESTING THE FIT OF LINEAR MODELS

SPOKOINY, V.

*Weierstrass Institute for Applied Analysis and Stochastics,
Mohrenstr. 39, 10117 Berlin*

1991 *Mathematics Subject Classification.* 62H25; Secondary 62G10.

Key words and phrases. data-driven test, Haar basis, linear hypothesis, nonparametric alternative, regression model.

ABSTRACT. The paper is concerned with the problem of testing a linear hypothesis about regression function. We propose a new testing procedure based on the Haar transform which is adaptive to unknown smoothness properties of the underlying function. The results show rate optimality of this procedure under mild conditions on the model.

1. Introduction

Suppose we are given data $(X_i, Y_i), i = 1, \dots, n$, with $X_i \in \mathbb{R}^1$, $Y_i \in \mathbb{R}^1$, obeying the regression equation

$$Y_i = f(X_i) + \xi_i \quad (1.1)$$

where f is an unknown regression function and ξ_i are zero mean random errors. Statistical analysis for such models may focus on the qualitative features of the underlying function f . Particularly, no-response model corresponds to testing the simple zero hypothesis that f is a constant function. Another typical example is connected to the hypothesis of linearity. More generally one may consider a parametric type hypothesis about f . In this paper, we restrict ourselves to the case of the hypothesis of linearity. Using the hypothesis testing framework, we test the null hypothesis $H_0 : f$ ‘is linear’, that is, $f(x) = a + bx$ for some constants a, b , versus the alternative $H_1 : f$ ‘is not linear’.

The problem of testing a simple or parametrically specified hypothesis is one of the classical in statistical inference, see e.g. Neyman (1937), Mann and Wald (1942), Lehmann (1957). Let ϕ be a test i.e. a measurable function of the observations Y_1, \dots, Y_n with two values 0, 1. As usual, the event $\{\phi = 0\}$ is treated as accepting the hypothesis and $\phi = 1$ means that the hypothesis is rejected. The quality of a test ϕ is described in terms of the corresponding error probabilities of the first and second kinds. Let \mathbf{P}_f denote the distribution of the data Y_1, \dots, Y_n for a fixed model function f , see (1.1). If f coincides with a linear function f_0 , then the error probability of the first kind at the point f_0 is the probability under f_0 to reject the hypothesis,

$$\alpha_{f_0}(\phi) = \mathbf{P}_{f_0}(\phi = 1).$$

Similarly one defines the error probability $\beta_f(\phi)$ of the second kind. If the function f is not linear, then

$$\beta_f(\phi) = \mathbf{P}_f(\phi = 0).$$

Typically one aims to construct a test φ of the prescribed level α_0 , that is, satisfying for a given $\alpha_0 > 0$ the condition $\alpha_{f_0}(\varphi) \leq \alpha_0$ which also has a nontrivial power $1 - \beta_f(\varphi) > 0$ against a possibly large class of alternatives f . A large number of proposals for

constructing such tests can be found in the literature. We refer to Hart (1997) where the reader can find historical remarks and further references. Note meanwhile, that the majority of results in this domain is concentrated either only on verifying the condition $\alpha_{f_0}(\phi) \leq \alpha_0$ or on studying asymptotic properties of the power function $1 - \beta_f(\varphi)$ for a fixed or local alternative, and the question of test optimality is not addressed rigorously.

One possibility to introduce test optimality is proposed by Ingster (1982). The idea is to construct a test ϕ which fulfills the above constraints $\alpha_{f_0}(\varphi) \leq \alpha_0$ for all linear functions f_0 and additionally the condition $\beta_f(\varphi) \leq \beta_0$ with some $\beta_0 < 1 - \alpha_0$ uniformly over a possibly large class \mathcal{F} of alternatives f . Following to Ingster (1982, 1993), we consider the class $\mathcal{F}(\varrho)$ consisting of smooth (in some sense) alternatives which are also separated from the set of linear functions with the distance ϱ , that is,

$$\inf_{a,b} \|f(\cdot) - a - b \cdot\| \geq \varrho,$$

$\|\cdot\|$ being the usual L_2 -norm. Then the quality of a test φ with the level α_0 is measured by a minimal distance ρ such that $\beta_f(\phi) \leq \beta_0$ for all f from $\mathcal{F}(\rho)$. A test ϕ^* with the level α_0 is *optimal* if it minimizes the corresponding separation distance ρ . Under this approach, the goal is both to evaluate the minimal possible separation distance ρ and to describe the corresponding optimal tests.

It turns out that the structure of optimal tests and the corresponding separation distance strongly depend on the smoothness class \mathcal{F} we consider. Ingster (1982, 1993) described the optimal rate of decay of the separation distance ρ to zero as the sample size n tends to infinity for Hölder and Sobolev function classes, the case of Besov classes is considered in Lepski and Spokoiny (1998). Sharp optimal asymptotic results can be found in Ermakov (1990), Lepski (1993), Lepski and Tsybakov (1996), Ingster and Suslina (1998).

Unfortunately the mentioned procedures hardly apply in practice since the information about smoothness properties of the underlying function f is typically lacking. Some adaptive (data-driven) smooth tests are proposed in Ledwina (1994), Fan (1996), Ledwina and Kallenberg (1997), Hart (1997) where the reader can find further references. Spokoiny (1996, 1998) considered the problem of adaptive testing against a smooth alternative and constructed an adaptive test which is near optimal by a *log log* multiple for a wide range of smoothness classes. Moreover, the test is rate optimal in the class of adaptive tests, that is, this *log log* factor is an unavoidable payment for the adaptive property. The inconvenience for practical applications is that this procedure is designed for an idealized ‘signal + white noise’ model and only the case of a simple null is considered.

The aim of this paper is to develop an adaptive testing method which allows for a non-regular design, non-Gaussian errors with an unknown distribution and a non-simple null, and which is computationally simple and stable w.r.t. the design non-regularity. The latter property is achieved by making use of the simplest wavelet basis, namely the

Haar transform. It is worth mentioning that the Haar basis is not often used for estimating the regression function f from (1.1) because of its non-regularity: the corresponding estimator is only rate suboptimal. Nevertheless, Ingster (1993) shown that, in spite of the non-regularity of the Haar basis, the corresponding testing procedure is rate optimal. Another remark concerns the assumption on the errors ξ_i . Assuming i.i.d. errors with a known distribution, one can easily select a critical level for any test statistic using the Monte-Carlo or other resampling technique. For practical applications, this approach needs to be justified since the underlying error distribution is typically unknown. The problem becomes even more complicated if a data-driven test basing on the maximum of different test statistics is used. We establish some general results on the approximation of quadratic forms of independent random variables by similar quadratic forms of Gaussian random variables which help to justify the following recipe: if the critical level of the considered test statistic is calculated for Gaussian errors, then it applies, at least asymptotically, as the sample size grows, for an arbitrary errors distribution with bounded 4 moments.

The paper is organized as follows. Section 2 contains the description of the proposed testing procedure. The properties of this procedure are discussed in Section 3. Some possible extensions of the method to the multivariate regression and heterogeneous noise can be found in Section 4. The proofs are postponed to Section 5. In the Appendix we collect some general results for quadratic forms.

2. Testing procedure

We consider the univariate regression model

$$Y_i = f(X_i) + \xi_i, \quad i = 1, \dots, n,$$

with additive homogeneous noise, that is, the errors ξ_i are independent identically distributed with zero mean and the variance σ^2 : $\mathbf{E}\xi_i = 0$ and $\mathbf{E}\xi_i^2 = \sigma^2$. The design points X_1, \dots, X_n are assumed to be rescaled to the interval $[0, 1]$, that is, $X_i \in [0, 1]$ for all $i = 1, \dots, n$.

The proposed makes use of the Haar transform. We first recall some useful facts about the Haar decomposition and then explain the idea of the method.

2.1. Preliminaries

Hereafter we denote by I the multi-index $I = (j, k)$ with $j = 0, 1, 2, \dots$ and $k = 0, 1, \dots, 2^j - 1$, and let \mathcal{I} be the set of all such multi-indices. We also set

$$\mathcal{I}_j = \{(j, k), k = 0, 1, \dots, 2^j - 1\}$$

for the index set corresponding to j -th level. Let now the function $\psi(t)$ be defined by

$$\psi(t) = \begin{cases} 0 & t < 0, \\ 1 & 0 \leq t < 1/2, \\ -1 & 1/2 \leq t < 1, \\ 0 & t > 1. \end{cases}$$

For every $I = (j, k)$, define the Haar basis function h_I by

$$h_I(t) = 2^{j/2} \psi(2^j t - k).$$

Clearly the function h_I is supported on the interval $A_I = [2^{-j}k, 2^{-j}(k+1)]$. It is well known that each measurable function f on $[0, 1]$ can be decomposed in the following way

$$f(t) = c_0 + \sum_{I \in \mathcal{I}} c_I h_I(t) = c_0 + \sum_{j=0}^{\infty} \sum_{I \in \mathcal{I}_j} c_I h_I(t). \quad (2.1)$$

This means that the problem of recovering the function f can be transformed into the problem of estimating the coefficients c_I by given data. Since we have only n observations, it makes no sense to estimate more (in order) than n coefficients. We restrict therefore the total number of considered levels j . Let some j be fixed such that $2^{j+1} \leq n$. We also introduce the rescaled basis functions ψ_I to provide $\sum_i |\psi_I(X_i)|^2 = 1$, that is,

$$\psi_I(X_i) = \mu_I^{-1} h_I(X_i),$$

with $\mu_I^2 = \sum_{i=1}^n h_I^2(X_i)$. Next we replace the infinite decomposition (2.1) by the finite approximation $\sum_{I \in \mathcal{I}(j)} c_I \psi_I(t)$ where the index set $\mathcal{I}(j)$ contains all level sets \mathcal{I}_ℓ with $\ell \leq j$. Taking into account the structure of the null hypothesis, we complement the set of functions $(\psi_I, I \in \mathcal{I}_\ell)$, $\ell \leq j$, with two functions $\psi_0 \equiv 1$ and $\psi_1(t) = t$, that is, we consider the set of indices

$$\mathcal{I}(j) = \{0, 1\} + \bigcup_{\ell=0}^j \mathcal{I}_\ell.$$

The idea of the proposed procedure is to estimate all the coefficients $(c_I, I \in \mathcal{I}(j))$ from the data Y_1, \dots, Y_n and then test that all the coefficients c_I for $I \neq 0, 1$ are zero.

For a function g , define $\|g\|_n$ by

$$\|g\|_n^2 = \frac{1}{n} \sum_{i=1}^n g^2(X_i).$$

Define also the column-vector $\boldsymbol{\theta}^*(j) = (\theta_I^*, I \in \mathcal{I}(j))$ as a minimizer of the error of approximating f by a linear combination of ψ_I , $I \in \mathcal{I}(j)$:

$$\boldsymbol{\theta}^*(j) = \underset{\boldsymbol{\theta}(j)}{\operatorname{arginf}} \|f - \sum_{I \in \mathcal{I}(j)} \theta_I \psi_I\|_n^2. \quad (2.2)$$

This is a quadratic optimization problem with respect to the coefficients $\{\theta_I, I \in \mathcal{I}(j)\}$. Therefore, the solution $\boldsymbol{\theta}^*$ always exists but it is probably non unique. To get an explicit representation for $\boldsymbol{\theta}^*$ we introduce matrix notation.

First of all, we make an agreement to identify every function g with the vector $(g(X_i), i = 1, \dots, n)^\top$ in \mathbb{R}^n where the symbol $^\top$ means transposition. Particularly, the model function f is identified with the vector $(f(X_i), i = 1, \dots, n)^\top$.

Denote by N_j the number of elements at each level j ,

$$N_j = \#(\mathcal{I}_j) = 2^j, \quad j = 0, 1, \dots, j$$

and let $N(j)$ be the total number of elements in the set $\mathcal{I}(j)$,

$$N(j) = 2 + \sum_{\ell=0}^j N_\ell = 1 + 2^{j+1}.$$

Introduce $n \times N(j)$ -matrix $\Psi(j) = (\psi_{i,I}, i = 1, \dots, n, I \in \mathcal{I}(j))$ with elements

$$\psi_{i,I} = \psi_I(X_i) = \psi_I(X_i), \quad I \in \mathcal{I}(j), i = 1, \dots, n.$$

Clearly $\psi_I(X_i) = \pm 1/\sqrt{M_I}$ where M_I is the number of design points in the interval A_I corresponding to the index I , and also $\psi_{i,0} = n^{-1/2}$ and $\psi_{i,1} = X_i (\sum_{\ell=1}^n X_\ell^2)^{-1/2}$. Now the approximation problem (2.2) can be rewritten in the form

$$\boldsymbol{\theta}^*(j) = \underset{\boldsymbol{\theta}(j)}{\operatorname{arginf}} \|f - \Psi(j)\boldsymbol{\theta}(j)\|_n^2.$$

The solution to this quadratic problem can be represented as

$$\boldsymbol{\theta}^*(j) = \left(\Psi(j)^\top \Psi(j) \right)^{-1} \Psi(j)^\top f. \quad (2.3)$$

Strictly speaking, this representation is valid only if the matrix $\Psi(j)^\top \Psi(j)$ is not degenerate. In the general case, one may use the similar expression for $\boldsymbol{\theta}^*(j)$ when understanding $(\Psi(j)^\top \Psi(j))^{-1}$ as a pseudo-inverse matrix.

If the function f is linear, that is, $f(x) = \theta_0 + \theta_1 x$, we clearly get $\theta_0^* = \theta_0$, $\theta_1^* = \theta_1$ and $\theta_I^* = 0$ for all $I = (\ell, k)$ with $\ell \geq 0$ and $k \geq 0$. For a non-linear function f , the sum $\sum_{\ell=0}^j \sum_{I \in \mathcal{I}_\ell} |\theta_I^*|^2$ can be used to characterize the deviation of f from the space of linear functions.

Since the function f is observed with a noise, we cannot calculate directly the coefficients θ_I^* and we consider the least squares estimator $\widehat{\boldsymbol{\theta}}(j)$ of the vector $\boldsymbol{\theta}^*(j)$ which is

defined by minimization of the sum of residuals squared,

$$\widehat{\boldsymbol{\theta}}(j) = \underset{\boldsymbol{\theta}(j)}{\operatorname{arginf}} \|Y - \Psi(j)\boldsymbol{\theta}(j)\|_n^2 = \underset{\{\theta_I \in \mathcal{I}(j)\}}{\operatorname{arginf}} \sum_{i=1}^n \left(Y_i - \sum_{I \in \mathcal{I}(j)} \theta_I \psi_I(X_i) \right)^2. \quad (2.4)$$

Here \mathbf{Y} means the column-vector with elements Y_i , $i = 1, \dots, n$.

Define $V(j)$ as the pseudo-inverse of $\Psi(j)^\top \Psi(j)$, $V(j) = (\Psi(j)^\top \Psi(j))^-$. It is a symmetric $N(j) \times N(j)$ matrix (by $v_{I,I'}$ we denote its elements, $I, I' \in \mathcal{I}(j)$) and

$$\widehat{\boldsymbol{\theta}}(j) = V(j)\Psi(j)^\top \mathbf{Y}. \quad (2.5)$$

Neyman (1937) proposed a ‘smooth’ test based on the centralized and standardized sum of squares $\sum_{\ell=0}^j \sum_{I \in \mathcal{I}_\ell} |\widehat{\theta}_I|^2$ for some j . Ingster (1982, 1993) suggested the special choice of j depending on the smoothness properties of the function f which allows for a rate optimal testing. We follow Spokoiny (1996) where the method of Ingster (1993) is extended to adaptive testing by considering all such tests for different j simultaneously. Here we slightly modify that approach and consider the family of levelwise tests, that is, for every level j , we construct a test statistic based only on the empirical Haar coefficients $\widehat{\theta}_I$ for $I \in \mathcal{I}_j$, and the resulting test is defined as the maximum of all levelwise ones.

Let some number $j(n)$ be fixed such that $2^{j(n)} \leq n$ and let, for every $j \leq j(n)$, the estimate $\widehat{\boldsymbol{\theta}}(j)$ be defined by (2.4). Denote by $\widehat{\boldsymbol{\theta}}_j$ the part of the vector $\widehat{\boldsymbol{\theta}}(j)$ corresponding to the level j ,

$$\widehat{\boldsymbol{\theta}}_j = (\widehat{\theta}_I, I \in \mathcal{I}_j).$$

We analyze every such vector separately for all $j \leq j(n)$. Namely, for every $j \leq j(n)$, we use the statistic based on the sum $\sum_{I \in \mathcal{I}_j} |\widehat{\theta}_I|^2$ corresponding to j th resolution level.

To define our test, we need to have a more detailed insight into the properties of such sums under the null hypothesis, i.e. when the function f is linear: $f(x) = \theta_0 + \theta_1 x$. We have already mentioned that in this situation $f = \Psi(j)\boldsymbol{\theta}^*$ where $\theta_0^* = \theta_0$, $\theta_1^* = \theta_1$ and all remaining coefficients θ_I^* vanish. Therefore, using the model equation $\mathbf{Y} = f + \boldsymbol{\xi}$, we obtain

$$\begin{aligned} \widehat{\boldsymbol{\theta}}(j) &= V(j)\Psi(j)^\top (f + \boldsymbol{\xi}) \\ &= V(j)\Psi(j)^\top \Psi(j)\boldsymbol{\theta}^* + V(j)\Psi(j)^\top \boldsymbol{\xi} \\ &= \boldsymbol{\theta}^* + V(j)\Psi(j)^\top \boldsymbol{\xi}. \end{aligned} \quad (2.6)$$

Obviously $\boldsymbol{\zeta}(j) = V(j)\Psi(j)^\top \boldsymbol{\xi}$ is a random vector in $\mathbb{R}^{N(j)}$ with zero mean. Moreover, it holds for its covariance matrix

$$\begin{aligned} \mathbf{E}\boldsymbol{\zeta}(j)\boldsymbol{\zeta}(j)^\top &= V(j)\Psi(j)^\top \mathbf{E}\boldsymbol{\xi}\boldsymbol{\xi}^\top \Psi(j)V(j) \\ &= \sigma^2 V(j)\Psi(j)^\top \Psi(j)V(j) \\ &= \sigma^2 V(j). \end{aligned} \quad (2.7)$$

Due to (2.6), the subvector $\widehat{\boldsymbol{\theta}}_j$ of $\widehat{\boldsymbol{\theta}}(j)$ coincides under the null with the corresponding subvector $\boldsymbol{\zeta}_j$ of the vector $\boldsymbol{\zeta}(j)$, and it holds under the null in view of (2.7)

$$\begin{aligned}\mathbf{E}\widehat{\boldsymbol{\theta}}_j &= \mathbf{E}\boldsymbol{\zeta}_j = 0, \\ \mathbf{E}\widehat{\boldsymbol{\theta}}_j\widehat{\boldsymbol{\theta}}_j^\top &= \mathbf{E}\boldsymbol{\zeta}_j\boldsymbol{\zeta}_j^\top = \sigma^2 V_j.\end{aligned}$$

This particularly implies

$$\mathbf{E} \sum_{I \in \mathcal{I}_j} |\widehat{\theta}_I|^2 = \mathbf{E} \sum_{I \in \mathcal{I}_j} |\zeta_I|^2 = \sigma^2 \operatorname{tr} V_j$$

where $\operatorname{tr} A$ denotes the trace of a matrix A . Moreover, for the case of Gaussian errors ξ_i in (1.1), the estimates $\widehat{\theta}_I$ are also Gaussian random variables, and it holds

$$\begin{aligned}\operatorname{Var} \left(\sum_{I \in \mathcal{I}_j} |\widehat{\theta}_I|^2 \right) &= \mathbf{E} \left(\sum_{I \in \mathcal{I}_j} |\widehat{\theta}_I|^2 - \sigma^2 \operatorname{tr} V_j \right)^2 \\ &= \mathbf{E} \left(\sum_{I \in \mathcal{I}_j} |\zeta_I|^2 - \sigma^2 \operatorname{tr} V_j \right)^2 = 2\sigma^4 \operatorname{tr} V_j^2,\end{aligned}\tag{2.8}$$

see (2.7). This leads to the obvious idea to use the centralized and normalized sum

$$T_j = \frac{1}{\sqrt{2\sigma^4 \operatorname{tr} V_j^2}} \left(\sum_{I \in \mathcal{I}_j} |\widehat{\theta}_I|^2 - \sigma^2 \operatorname{tr} V_j \right)$$

as a test statistic. To define our testing procedure, we simply take the maximum of all such statistics over the set of all considered Haar levels j .

2.2. Testing procedure

First we define the finest considered resolution level $j(n)$ which has to satisfy $n2^{j(n)} \rightarrow \infty$, e.g.

$$j(n) = [\log_2 n - \log_2 \log_2 n].$$

where $[a]$ denotes the integer part of a . For each $j \leq j(n)$, let $\widehat{\boldsymbol{\theta}}(j)$ be defined by (2.5). Denote by $\widehat{\boldsymbol{\theta}}_j$ the part of the vector $\widehat{\boldsymbol{\theta}}(j)$ corresponding to the level j ,

$$\widehat{\boldsymbol{\theta}}_j = (\widehat{\theta}_I, I \in \mathcal{I}_j)$$

and let V_j be the submatrix of the matrix $V(j) = (\Psi(j)^\top \Psi(j))^{-1}$ corresponding to the level j , i.e. $V_j = (v_{I,I'}, I, I' \in \mathcal{I}_j)$. We consider χ^2 -type statistics

$$S_j = \|\widehat{\boldsymbol{\theta}}_j\|^2 = \sum_{I \in \mathcal{I}_j} \widehat{\theta}_I^2.$$

and define test statistics T_j by centralization and normalization of S_j :

$$T_j = \frac{1}{\sqrt{2\hat{\sigma}^4 \operatorname{tr} V_j^2}} \left(\sum_{I \in \mathcal{I}_j} |\hat{\theta}_I|^2 - \hat{\sigma}^2 \operatorname{tr} V_j \right)$$

where $\hat{\sigma}$ is the estimate of the error standard deviation defined in the next subsection. The proposed test rejects the null hypothesis, if at least one such statistic is significantly large, that is,

$$\phi^* = \mathbf{1}(T^* > \lambda) \quad \text{with} \quad T^* = \max_{j=0, \dots, j(n)} |T_j|$$

where λ is a critical value. The choice of λ is discussed in Section 2.4.

2.3. Estimation of σ^2

Recall that we assume a homogeneous additive noise in the model (1.1), that is, the errors ξ_i are independent identically distributed random variables fulfilling $\mathbf{E}\xi_i = 0$ and $\mathbf{E}\xi_i^2 = \sigma^2$. The variance σ^2 is typically unknown in practical applications but this value is important for the definition of our test procedure. Below we discuss how it can be estimated from the data Y_1, \dots, Y_n . We suppose for simplicity that the design points are ordered in a way that $X_1 \leq \dots \leq X_n$. There are several proposals for variance estimation. One possibility is to estimate σ^2 by the expression of the form $\frac{1}{2(n-1)} \sum_{i=1}^{n-1} (Y_{i+1} - Y_i)^2$, see Gasser et al. (1986). We follow the proposal from Hart (1997, Section 5.3) which provides an unbiased estimate of the variance under the linear null hypothesis.

Define *pseudo-residuals*

$$\begin{aligned} \hat{e}_i &= \frac{(X_{i+1} - X_i)}{(X_{i+1} - X_{i-1})} Y_{i-1} + \frac{(X_i - X_{i-1})}{(X_{i+1} - X_{i-1})} Y_{i+1} - Y_i \\ &= a_i Y_{i-1} + b_i Y_{i+1} - Y_i, \quad i = 2, \dots, n-1. \end{aligned}$$

which are the result of joining Y_{i+1} and Y_{i-1} by a straight line and taking the difference between this line and Y_i . A variance estimate based on these pseudo-residuals is

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=2}^{n-1} \frac{\hat{e}_i^2}{a_i^2 + b_i^2 + 1}. \quad (2.9)$$

It is obvious that $\mathbf{E}\hat{\sigma}^2 = \sigma^2$ if f is a linear function. Some other properties of this estimates are listed in Lemmas 5.1, 5.2 and 5.9 below.

2.4. Critical value λ

Here we discuss how to select the critical value λ to provide, at least asymptotically for large n , the condition $\alpha_{f_0}(\phi^*) \leq \alpha_0$ for all linear functions f_0 . We apply a bootstrap

procedure resampling from the no-response model (which is a particular case of a linear model) with standard normal errors

$$Y_{i,m}^* = \xi_{i,m}^*, \quad i = 1, \dots, n,$$

for $m = 1, \dots, M$, where the design points X_1, \dots, X_n are the same as for the original model (1.1), ξ_1^*, \dots, ξ_n^* are i.i.d. standard normal random variables and M is the considered number of bootstrap samples.

For every bootstrap sample $Y_{1,m}^*, \dots, Y_{n,m}^*$, we recalculate the test statistic T_m^* from this sample using the previous procedure (including the step of variance estimation). Finally we define the critical value λ as the α_0 -level for the set $\{T_m^*, m = 1, \dots, M\}$:

$$\lambda = \min \left\{ t : M^{-1} \sum_{m=1}^M \mathbf{1}(T_m^* > t) \leq \alpha_0 \right\}.$$

3. Main results

In this section we present the results describing asymptotic properties of the proposed testing procedure. We first discuss the properties of the test under the null and then we consider the power of the test.

3.1. Behavior under the null

Let ϕ^* be the test introduced above. Our first result concerns with the case of Gaussian errors ξ_i in the model (1.1). In this situation, independently of the design, the nominal level of the test ϕ^* is exactly α_0 .

Theorem 3.1. *Let observations Y_i, X_i , $i = 1, \dots, n$, obey the regression model (1.1) with a deterministic design X_1, \dots, X_n and with i.i.d Gaussian errors $\xi_i \sim \mathcal{N}(0, \sigma^2)$. If the function f is linear, $f(x) = \theta_0 + \theta_1 x$, then*

$$\alpha_f(\phi^*) \equiv \mathbf{P}_f(\phi^* = 1) = \alpha_0.$$

Our next result deals with a more general situation when the errors ξ_i are i.i.d. with 6 finite moments. In this case we also need some mild regularity conditions on the design.

Recall the notation $A_I = [2^{-j}k, 2^{-j}(k+1)]$ and let M_I stand for the number of design points in A_I : $M_I = \#\{i : X_i \in A_I\}$. Design regularity particularly means that each interval A_I contains enough design points X_i .

(D) (i) It holds for some positive constants C_* and C^* and all $j \leq j(n)$

$$\begin{aligned} \inf_{I \in \mathcal{I}_j} 2^j M_I / n &\geq C_*, \\ \sup_{I \in \mathcal{I}_j} 2^j M_I / n &\leq C^*; \end{aligned}$$

(ii) For some fixed constant C_D and all $j \leq j(n)$

$$\text{tr } V_j^2 \geq C_D 2^j;$$

(iii) For some fixed constant C_V and all $j \leq j(n)$

$$\|V(j)\| \leq C_V.$$

Here the norm $\|A\|$ of a symmetric matrix A is understood as the maximal eigenvalue of this matrix.

Condition (D) is trivially fulfilled with $C_* = C^* = C_D = C_V = 1$ for the case of a uniform random or deterministic equidistant design when $V(j)$ is the unit matrix.

Theorem 3.2. *Let observations Y_i, X_i , $i = 1, \dots, n$, obey the regression model (1.1) with a deterministic design X_1, \dots, X_n satisfying (D) and with i.i.d. errors ξ_i satisfying $\mathbf{E}\xi_i = 0$, $\mathbf{E}\xi_i^2 = \sigma^2$ and $\mathbf{E}|\xi_i^2 - \sigma^2|^3 \leq \sigma^6 C_6$ where C_6 is a fixed constant. If the function f is linear, $f(x) = \theta_0 + \theta_1 x$, then*

$$\alpha_f(\phi^*) \equiv \mathbf{P}_f(\phi^* = 1) \leq \alpha_0 + \delta_1(n),$$

where $\delta_1(n)$ depends on n , C_6 and the constants C_*, C^*, C_D, C_V from condition (D) only and $\delta_1(n) \rightarrow 0$ as $n \rightarrow \infty$.

3.2. Sensitivity of the test

Now we state the results concerning the sensitivity of the proposed test ϕ^* . The first assertion presents sufficient conditions for detecting an alternative with a high probability. Next we demonstrate how these conditions can be transferred into a more usual form about the rate of testing against a smooth alternative.

Proposition 3.1. *Let the design X_1, \dots, X_n obey (D) and the errors ξ_1, \dots, ξ_n fulfill the conditions of Theorem 3.2. Let then the regression function f be differentiable with the Lipschitz continuous first derivative f' :*

$$|f'(s) - f'(t)| \leq L|s - t| \tag{3.1}$$

with some fixed constant L . Let also $\boldsymbol{\theta}_j^* = (\theta_I^*, I \in \mathcal{I}_j)$ be the subvector of the vector $\boldsymbol{\theta}^*(j)$ from (2.3) corresponding to j th resolution level and let V_j be the corresponding covariance submatrix, $j = 1, \dots, j(n)$. If, for some $j \leq j(n)$, it holds

$$T_j^* \equiv \frac{\|\boldsymbol{\theta}_j^*\|^2}{\sigma^2 \sqrt{2 \text{tr } V_j^2}} \geq 3(\lambda_n^{1/2} + 1)^2,$$

with $\lambda_n = \max\{\lambda, 2\sqrt{\log j(n)}\}$, then

$$\mathbf{P}(\phi^*(j) = 0) \leq \delta(n) \rightarrow 0, \quad n \rightarrow \infty,$$

where $\delta(n)$ depends on n , L and the constants C_6, C_*, C^*, C_D, C_V only.

We shall show, see Lemma 5.2 that, at least for sufficiently large n , it holds $\lambda \leq 2\sqrt{\log j(n)}(1 + o_n(1))$. Hence, the result of Proposition 3.1 means that the test ϕ^* detects with a probability close to one any alternative for which at least one from the corresponding values T_j^* exceeds $6\sqrt{\log j(n)}(1 + o_n(1))$. Therefore, the error of the second kind may occur with a significant probability only if

$$T_j^* \leq 6\sqrt{\log j(n)}(1 + o_n(1)), \quad 0 \leq j \leq j(n).$$

It remains to understand what follows for the function f from these inequalities.

3.3. Rate of testing against a smooth alternative

To formulate the results on the rate of testing, we have to introduce some smoothness conditions on the function f . This can be done in different ways. We choose one based on the accuracy of approximation of this function by piecewise polynomials of certain degree s . Given $j \leq j(n)$, denote by $\{A_I, I \in \mathcal{I}_j\}$ the partition of the interval $[0, 1]$ into intervals of length 2^{-j} : if $I = (j, k)$, then $A_I = [k2^{-j}, (k+1)2^{-j})$. Next, for an integer s , define $\mathcal{P}_s(j)$ as the set of piecewise polynomials of degree $s-1$ on the partition $\{A_I\}$ i.e. every function g from $\mathcal{P}_s(j)$ coincides on each A_I with a polynomial $a_0 + a_1x + \dots + a_{s-1}x^{s-1}$ where the coefficients a_0, \dots, a_{s-1} may depend on I . Now the condition that a function f has regularity s can be understood in the sense that this function is approximated by functions from $\mathcal{P}_s(j)$ at the rate 2^{-js} , or, more precisely,

$$\inf_{g \in \mathcal{P}_s(j)} \left[\int_0^1 |f(t) - g(t)|^2 dt \right]^{1/2} \leq C_s 2^{-js}$$

where a positive constant C_s depends on s only.

In our conditions we change the integral by summation over observation points. This helps to present the results in a more readable form without changing the sense of required conditions. It can be easily seen that if the design is regular, then the both forms are equivalent up to a constant factor.

Let now a function f be fixed. Let also j_0 be such that $2^{j_0-1} \geq s$. Set for $j \geq j_0$

$$r_s(j) = \inf_{g \in \mathcal{P}_s(j-j_0)} \|f - g\|_n = \inf_{g \in \mathcal{P}_s(j-j_0)} \left[\sum_{i=1}^n |f(X_i) - g(X_i)|^2 \right]^{1/2}.$$

The quantity $r_s(j)$ characterizes the accuracy of approximation of f by piecewise polynomials. In particular, the Haar approximation we use corresponds to the case when $s = 1$.

Theorem 3.3. *Let condition (D) hold, the errors ξ_1, \dots, ξ_n fulfill the conditions of Theorem 3.2, and the regression function f obey (3.1). There exist a constant \varkappa depending on the values C_V, C_D, C_*, C^* and L only, such that if f satisfies, for some*

$j \leq j(n)$, the following inequality

$$\inf_{a,b} \|f - a - b\psi_1\|_n \geq \varkappa \left(r_s(j) + \sqrt{2^{j/2} \lambda_n} \right) \quad (3.2)$$

with $\psi_1(x) = x$, then

$$\mathbf{P}_f(\phi^* = 0) \leq \delta(n) \rightarrow 0, \quad n \rightarrow \infty,$$

where $\delta(n)$ is shown in Proposition 3.1.

Remark 3.1. It is of interest to compare this result with existing results on the rate of hypothesis testing. For instance, it was shown in Ingster (1982) that if f belongs to a Sobolev ball $W_s(1)$ with

$$W_s(1) = \left\{ f : \int_0^1 |f^{(s)}(x)|^2 dx \leq 1 \right\},$$

$f^{(s)}$ being s th derivative of f , then the optimal rate of testing is $n^{-2s/(4s+1)}$.

For our procedure, the following result is a straightforward corollary of Theorem 3.3.

Corollary 3.1. *Let the underlying function f belong to a Sobolev ball $W_s(1)$ and let condition (D) hold. There exists a constant $C_s > 0$ depending on s and the constants from condition (D) only and such that, for n large enough, the inequality*

$$\inf_{a,b} \|f - a - b\psi_1\|_n^2 \geq C_s (n/\lambda_n)^{-\frac{2s}{4s+1}}$$

implies

$$\mathbf{P}(\phi^* = 0) = o_n(1).$$

We observe that the proposed method is rate near optimal by a log-log multiple.

Remark 3.2. The result of Theorem 3.3 helps to understand what happens in the case when the design is not regular and, for instance, if there some intervals I with $M_I = 0$. It was already mentioned that the procedure applies in this situation as well and the error probability of the first kind is about α_0 at least for n sufficiently large and for Gaussian errors ξ_i . Concerning the error probability of the second kind, the inspection of the proof shows that design irregularity decreases the sensitivity of our procedure in the following sense: there exist smooth alternatives with probably large L_2 -norm which are not detected. This may occur e.g. in the situation when f is deviated from the best linear approximation only in the domain with very few design points inside.

4. Some extension

Here we briefly discuss some possible extensions of the procedure.

4.1. Heterogeneous noise

The proposed procedure essentially uses the noise homoskedasticity. Namely, this condition allows to estimate the unknown noise variance at the rate $n^{-1/2}$. If this assumption is not fulfilled (this is for instance the case for the binary response model, see e.g. Klein and Spady, 1993), then the direct application of the method from Section 2 becomes questionable. One often used approach in such situation is based on some local estimation of the variance as a function of the design point x . Unfortunately, this may lead to a very poor quality of variance estimation for small and moderate sample size n . This may in turn destroy the behaviour of the test both under the null and the alternative because the variance estimate is used for centering the considered test statistics. A more useful approach is to avoid centering either by splitting the sample into two independent subsamples or by removing the diagonal terms from the considered test statistics. We briefly discuss the latter possibility. Each empirical Haar coefficient $\widehat{\theta}_I$ is a linear combination of the observations Y_i , see (2.5). Denote by $w_{I,i}$ the corresponding coefficients: $\widehat{\theta}_I = \sum_{i=1}^n w_{I,i} Y_i$. Then clearly

$$|\widehat{\theta}_I|^2 = \sum_{i=1}^n \sum_{i'=1}^n w_{I,i} w_{I,i'} Y_i Y_{i'}.$$

To define our modified test statistics, we remove from this sum the diagonal elements with $i = i'$:

$$T'_j = \frac{\sum_{I \in \mathcal{I}_j} |\widehat{\theta}_I|^2 - \sum_{I \in \mathcal{I}_j} \sum_{i=1}^n w_{I,i}^2 Y_i^2}{\sqrt{2 \operatorname{tr} V_j^2}}.$$

The critical level for the test statistic $T^* = \max_{j \leq j(n)} \{T'_j\}$ can be again calculated by the bootstrap procedure when resampling from the heterogeneous model with $Y_i^* = \sigma_i \xi_i^*$ where $\sigma_i^2 = (a_i^2 + b_i^2 + 1)^{-1} \widehat{e}_i^2$ and the pseudo-residuals \widehat{e}_i are defined in Section 2.3.

4.2. Linear parametric hypothesis

The proposed method allows for the straightforward generalization to the case of a linear null of the form $f(x) = \theta_1 \psi_1(x) + \dots + \theta_p \psi_p(x)$ with known function g_1, \dots, g_p . One should simply include this function in the set $\{\psi_I, I \in \mathcal{I}(j)\}$ and then proceed as before. For theoretical study, the only properties of the estimate $\widehat{\sigma}^2$ of σ^2 have to be refined.

4.3. General parametric hypothesis

The situation becomes more complicated for a general parametric null. Here one possibility is, similarly to Härdle and Mammen (1993), to construct first the parametric fit, then to subtract it from the data and finally to apply the above procedure for testing a no-response hypothesis. A more detailed study of such test needs to be done.

4.4. Multivariate regression

The procedure allows also for straightforward generalization to the multivariate regression case. We may use the corresponding multivariate Haar basis taking so many levels that the total number of estimated coefficients does not exceed n . Some further extensions to additive or generalized additive models are also possible, see e.g. Härdle et al. (1998).

5. Proofs

In this section we first prove Theorems 3.1 and 3.3 for the case of Gaussian errors ξ_i and then discuss the generalization to the general case.

5.1. Proof of Theorem 3.1

It suffices to check that the distribution of the test statistic T^* based on the bootstrap sample Y_1^*, \dots, Y_n^* is the same as for the original sample Y_1, \dots, Y_n . The difference between these two samples is only in the linear trend (which can be nontrivial for the original sample but does not appear in the bootstrap one) and in the noise variance (we resample with the error variance 1 instead of σ^2). Note however that the linear trend in the regression function makes no influence on the considered test statistics T_j . Indeed, the numerator of this statistic is defined as the centered sum over \mathcal{I}_j of the the empirical Haar coefficients $\hat{\theta}_I$ squared, so that the coefficients $\hat{\theta}_0$ and $\hat{\theta}_1$, corresponding to the linear trend, do not enter, see (2.7) and (2.6). Similarly, the estimate $\hat{\sigma}^2$ of the noise variance σ^2 is based on the pseudo-residuals \hat{e}_i which are defined in a way that the linear trend in the regression function cancels out, see Lemma 5.1.

Further, for the case of zero trend, both numerator and denominator of each T_j is some quadratic forms of the errors ξ_i which can be represented as $\xi_i = \sigma \tilde{\xi}_i$ with i.i.d. standard normal variables $\tilde{\xi}_i$, $i = 1, \dots, n$. This yields, see (2.9), that the distribution of each test statistic T_j does not depend on σ . The same is obviously true for the maximum T^* and the assertion follows.

5.2. The properties of the estimate $\hat{\sigma}^2$

Here we discuss the properties of the estimate $\hat{\sigma}^2$ of the noise variance σ^2 . We present two results. The first one describes the properties under the null, and the second one applies under a smooth alternative as well. The results are stated under the Gaussian errors ξ_i . For the extension, see Section 5.5.

Lemma 5.1. *Let the regression function f be linear. Then for $n \geq 36$ and each $\gamma \geq 1$ with $\gamma \leq \frac{3}{7} \sqrt{\frac{3(n-2)}{2}}$*

$$\mathbf{P} \left(\pm \sqrt{n} \left(\frac{\hat{\sigma}^2}{\sigma^2} - 1 \right) > 2\gamma \right) \leq e^{-\gamma^2/4},$$

Proof. For the case of a linear function $f(x) = \theta_0 + \theta_1 x$, one easily gets with the coefficients $a_i = \frac{(X_{i+1} - X_i)}{(X_{i+1} - X_{i-1})}$, $b_i = \frac{(X_i - X_{i-1})}{(X_{i+1} - X_{i-1})}$

$$a_i f(X_{i-1}) + b_i f(X_{i+1}) - f(X_i) = 0.$$

Now the model equation (1.1) implies

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=2}^{n-1} |\eta_i|^2$$

with

$$\eta_i = \frac{a_i \xi_{i-1} + b_i \xi_{i+1} - \xi_i}{\sqrt{a_i^2 + b_i^2 + 1}}.$$

To estimate the difference $|\hat{\sigma}^2 - \sigma^2|$, we apply Proposition 6.1. Let $\boldsymbol{\eta}$ denote the vector $(\eta_2, \dots, \eta_{n-1})^\top$. Obviously $\mathbf{E}\boldsymbol{\eta} = 0$. Define $\Sigma = \mathbf{E}\boldsymbol{\eta}\boldsymbol{\eta}^\top$. Observe first that

$$\frac{1}{n-2} \text{tr} \Sigma = \frac{1}{n-2} \sum_{i=2}^{n-1} \frac{\sigma^2 (a_i^2 + b_i^2 + 1)}{(a_i^2 + b_i^2 + 1)} = \sigma^2.$$

Next, it is easy to check that $1/2 \leq a_i^2 + b_i^2 \leq 1$ and

$$\frac{\sigma^2}{2} \leq \frac{\sigma^2 (a_i + b_i)}{\sqrt{(a_i^2 + b_i^2 + 1)(a_{i+1}^2 + b_{i+1}^2 + 1)}} \leq \frac{2\sigma^2}{3}.$$

Hence

$$\begin{aligned} \mathbf{E}\eta_i^2 &= \sigma^2, \\ |\mathbf{E}\eta_i \eta_{i+1}| &\leq 2\sigma^2/3, \\ |\mathbf{E}\eta_i \eta_{i+1}| &\geq \sigma^2/2, \\ \mathbf{E}\eta_i \eta_{i'} &= 0, \quad |i' - i| > 1, \end{aligned}$$

This allows to estimate $\text{tr } \Sigma^2$ as follows:

$$\begin{aligned}
\frac{1}{(n-2)^2} \text{tr } \Sigma^2 &= \frac{1}{(n-2)^2} \sum_{i=2}^{n-1} \sum_{j=2}^{n-1} (\mathbf{E} \eta_i \eta_j)^2 \\
&= \frac{1}{(n-2)^2} \sum_{i=2}^{n-1} \sum_{j=2}^{n-1} [(\mathbf{E} \eta_{i-1} \eta_i)^2 + (\mathbf{E} \eta_i^2)^2 + (\mathbf{E} \eta_i \eta_{i+1})^2] \\
&\leq \frac{1}{(n-2)^2} \sigma^4 \sum_{i=2}^{n-1} (1 + 4/9 + 4/9) \\
&= \frac{17\sigma^4}{9(n-2)} \leq \frac{2\sigma^4}{n}
\end{aligned}$$

for $n \geq 36$. Similarly

$$\begin{aligned}
\frac{1}{(n-2)^2} \text{tr } \Sigma^2 &\geq \frac{1}{(n-2)^2} \sigma^4 \sum_{i=2}^{n-1} (1 + 1/4 + 1/4) = \frac{3\sigma^4}{2(n-2)} \\
\frac{1}{n-2} \|\Sigma\| &\leq \frac{1}{n-2} \max_{i=1, \dots, n} \sum_{j=1}^n |\mathbf{E} \eta_i \eta_j| \leq \frac{7\sigma^2}{3(n-2)}.
\end{aligned}$$

This implies $\|\Sigma\|^{-1} \sqrt{\text{tr } \Sigma^2 / 2} \geq \frac{3}{7} \sqrt{\frac{3(n-2)}{2}}$ and the application of Lemma 6.1 with $\varepsilon = \frac{\eta}{\sqrt{n-2}}$ yields for every γ with $1 \leq \gamma \leq \frac{3}{7} \sqrt{\frac{3(n-2)}{2}}$

$$\mathbf{P} \left(\pm(\hat{\sigma}^2 - \sigma^2) > \gamma \sqrt{\frac{4\sigma^4}{n}} \right) \leq e^{-\gamma^2/4}$$

and the required assertion follows. □

Next we show that $\hat{\sigma}^2$ estimate the true value σ^2 at the rate $n^{-1/2}$ under a mild assumption on the regression function f and the design X_1, \dots, X_n . We again assume that the design points are renumbered to provide $X_1 \leq X_2 \leq \dots \leq X_n$.

Lemma 5.2. *Let the regression function f from (1.1) satisfies the condition*

$$|f'(s) - f'(t)| \leq L|s - t|$$

for some $L \geq 0$ and all s, t from $[0, 1]$. Let also the design X_1, \dots, X_n fulfill

$$X_{i+1} - X_i \leq Dn^{-1} \tag{5.1}$$

with some constant D . Then, for $n \geq 36$ and every γ with $1 \leq \gamma \leq \frac{3}{7} \sqrt{\frac{3(n-2)}{2}}$,

$$\mathbf{P} \left(\pm\sqrt{n} \left(\frac{\hat{\sigma}^2}{\sigma^2} - 1 \right) > 2\gamma(1 + n^{-1/2}) \right) \leq 2e^{-\gamma^2/4},$$

provided that

$$\frac{D^4 L^4 n^{-4}}{6} + \sqrt{\frac{D^4 L^4 n^{-4}}{6}} \left(\frac{4\sigma^4}{n} \right)^{1/4} < \frac{2\sigma^2}{n}. \tag{5.2}$$

Remark 5.1. The condition (5.2) is obviously fulfilled if $(DL)^2 < n\sigma$.

Proof. The definition of the coefficients a_i and b_i , see Section 2.3, provides for any linear function $\ell(x)$ the identity $a_i\ell(X_{i-1}) + b_i\ell(X_{i+1}) - \ell(X_i) = 0$. Now the smoothness properties of the function f imply for $\ell(x) = f(X_i) + f'(X_i)(x - X_i)$

$$|f(x) - \ell(x)| \leq 0.5L^2|x - X_i|^2$$

and hence, using (5.1) and the conditions $a_i \geq 0$, $b_i \geq 0$ and $a_i + b_i = 1$

$$\begin{aligned} & |a_i f(X_{i-1}) + b_i f(X_{i+1}) - f(X_i)| \\ &= |a_i [f(X_{i-1}) - \ell(X_{i-1})] + b_i [f(X_{i+1}) - \ell(X_{i+1})] - [f(X_i) - \ell(X_i)]| \\ &\leq 0.5L^2 a_i |X_i - X_{i-1}|^2 + 0.5L^2 b_i |X_{i+1} - X_i|^2 \\ &\leq 0.5D^2 L^2 n^{-2}. \end{aligned} \tag{5.3}$$

Next, define

$$\begin{aligned} \eta_i &= \frac{a_i \xi_{i-1} + b_i \xi_{i+1} - \xi_i}{\sqrt{a_i^2 + b_i^2 + 1}} \\ \Delta_i &= \frac{a_i f(X_{i-1}) + b_i f(X_{i+1}) - f(X_i)}{\sqrt{a_i^2 + b_i^2 + 1}}. \end{aligned}$$

Then

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=2}^{n-1} |\Delta_i + \eta_i|^2.$$

To estimate the difference $|\hat{\sigma}^2 - \sigma^2|$, we apply Proposition 6.2. Let $\eta = (\eta_2, \dots, \eta_{n-1})^\top$. We know, see the proof of Lemma 5.1, that $\mathbf{E}\eta = 0$ and the matrix $\Sigma = \mathbf{E}\eta\eta^\top$ fulfills

$$\begin{aligned} \frac{1}{n-2} \operatorname{tr} \Sigma &= \sigma^2, \\ \frac{1}{(n-2)^2} \operatorname{tr} \Sigma^2 &\leq \frac{2\sigma^4}{n}. \end{aligned}$$

The inequality $a_i^2 + b_i^2 \geq 1/2$ and (5.3) provide

$$\frac{1}{n-2} \|\Delta\|^2 = \frac{1}{n-2} \sum_{i=2}^{n-1} \Delta_i^2 \leq \frac{1}{n-2} \sum_{i=2}^{n-1} \frac{D^4 L^4 n^{-4}}{4(a_i^2 + b_i^2 + 1)} \leq \frac{D^4 L^4 n^{-4}}{6}.$$

The application of Proposition 6.2 with $\mathbf{c} = \frac{\Delta}{\sqrt{n-2}}$ and $\varepsilon = \frac{\eta}{\sqrt{n-2}}$ yields for every γ with $1 \leq \gamma \leq \frac{3}{7} \sqrt{\frac{3(n-2)}{2}}$

$$\mathbf{P} \left(\pm(\hat{\sigma}^2 - \sigma^2) > \frac{D^4 L^4 n^{-4}}{6} + \gamma \sqrt{\frac{D^4 L^4 n^{-4}}{6}} \left(\frac{4\sigma^4}{n} \right)^{1/4} + \gamma \sqrt{\frac{4\sigma^4}{n}} \right) \leq 2e^{-\gamma^2/4}$$

and the required assertion follows in view of (5.2). \square

Lemma 5.3. Let $N_j = 2^j$ denote the number of elements in the set \mathcal{I}_j . It holds

$$\frac{\operatorname{tr} V_j}{\sqrt{2 \operatorname{tr} V_j^2}} \leq \sqrt{N_j/2}.$$

Proof. Clearly

$$\operatorname{tr} V_j^2 = \sum_{I \in \mathcal{I}_j} \sum_{I' \in \mathcal{I}_j} v_{I,I'}^2 \geq \sum_{I \in \mathcal{I}_j} v_{I,I}^2.$$

Next, the Cauchy-Schwarz inequality implies

$$N_j^{-1} \operatorname{tr} V_j = N_j^{-1} \sum_{I \in \mathcal{I}_j} v_{I,I} \leq \left(N_j^{-1} \sum_{I \in \mathcal{I}_j} v_{I,I}^2 \right)^{1/2}$$

and the assertion follows. \square

Lemma 5.4. Let λ be the critical value of the test selected by the testing procedure. If design X_1, \dots, X_n fulfills (D), then, for n sufficiently large,

$$\lambda \leq 2\sqrt{\log j(n)}(1 + o_n(1)).$$

Proof. Recall that the critical value λ corresponds to the α_0 -value of the test statistic $T^* = \max_{j \leq j(n)} T_j$ under the no-response model $f(x) \equiv 0$ and under the assumption of standard normal errors ξ_i , $i = 1, \dots, n$. In such a situation, the subvector $\hat{\theta}_j$ of $\hat{\theta}(j)$ coincides with the Gaussian vector $\zeta_j \sim \mathcal{N}(0, V_j)$, see Section 2.1, and hence the corresponding statistic T_j can be represented in the form

$$T_j = \frac{\|\zeta_j\|^2 - \hat{\sigma}^2 \operatorname{tr} V_j}{\hat{\sigma}^2 \sqrt{2 \operatorname{tr} V_j^2}}.$$

and it suffices to show that

$$\mathbf{P} \left(\max_{j \leq j(n)} T_j > 2\sqrt{\log j(n)}(1 + \delta_1(n)) \right) \leq \delta_2(n)$$

with two numeric sequences $\delta_1(n) \rightarrow 0$ and $\delta_2(n) \rightarrow 0$.

Now, for every $z \geq 1$ and $a \in (0, 1)$,

$$\begin{aligned} \left\{ T_j > \frac{z+1}{a} \right\} &= \left\{ \frac{\|\zeta_j\|^2 - \hat{\sigma}^2 \operatorname{tr} V_j}{\sigma^2 \sqrt{2 \operatorname{tr} V_j^2}} > \frac{(z+1)\hat{\sigma}^2}{a\sigma^2} \right\} \\ &\subseteq \left\{ \frac{\|\zeta_j\|^2 - \sigma^2 \operatorname{tr} V_j}{\sigma^2 \sqrt{2 \operatorname{tr} V_j^2}} > z \right\} \cup \left\{ \frac{(\hat{\sigma}^2 - \sigma^2) \operatorname{tr} V_j}{\sigma^2 \sqrt{2 \operatorname{tr} V_j^2}} > 1 \right\} \cup \left\{ \frac{\hat{\sigma}^2}{\sigma^2} < a \right\}. \end{aligned}$$

This clearly yields in view of Lemma 5.3

$$\begin{aligned} & \mathbf{P} \left(\max_{j \leq j(n)} T_j > \frac{z+1}{a} \right) \\ & \leq \mathbf{P} \left(\frac{\widehat{\sigma}^2}{\sigma^2} < a \right) + \mathbf{P} \left(\frac{\widehat{\sigma}^2}{\sigma^2} - 1 > \frac{1}{\sqrt{N_{j(n)}/2}} \right) + \sum_{j=0}^{j(n)} \mathbf{P} \left(\frac{\|\zeta_j\|^2 - \sigma^2 \operatorname{tr} V_j}{\sigma^2 \sqrt{2 \operatorname{tr} V_j^2}} > z \right). \end{aligned}$$

We apply this bound with $z = 1 + v_n$ and $a = 1 - v_n^{-1}$ where $v_n = 2\sqrt{\log j(n)}$. It follows from condition (D) that $v_n \leq \|V_j\|^{-1} \sqrt{\operatorname{tr} V_j^2/2}$ for all $j \geq j_1$ where j_1 is the minimal integer satisfying $C_D 2^{j_1} > C_V^2 v_n^2$. An application of Proposition 6.1 with $\gamma = v_n$ and $t = 1$ for $j \geq j_1$ and with $\gamma = 1$ and $t = v_n$ allows to bound

$$\mathbf{P} \left(\frac{\|\zeta_j\|^2 - \sigma^2 \operatorname{tr} V_j}{\sigma^2 \sqrt{2 \operatorname{tr} V_j^2}} > v_n + 1 \right) \leq \begin{cases} e^{-v_n^2/4 - v_n/2} & j \geq j_1, \\ e^{-v_n/2} & \text{otherwise.} \end{cases}$$

Using also Lemma 5.1 we derive

$$\begin{aligned} & \mathbf{P} \left(T^* > \frac{2 + v_n}{1 - v_n^{-1}} \right) \\ & \leq e^{-n v_n^{-2}/4} + e^{-n/(2N_{j(n)})} + \sum_{j=0}^{j_1-1} e^{-v_n/2} + \sum_{j=j_1}^{j(n)} e^{-v_n^2/4 - v_n/2} \\ & \leq e^{-n v_n^{-2}/4} + e^{-n/(2N_{j(n)})} + \log_2(C_V^2 v_n^2 / C_D) e^{-v_n/2} + \frac{1+j(n)}{j(n)} e^{-v_n/2} \rightarrow 0, \quad n \rightarrow \infty. \end{aligned}$$

Here we have used that $n^{-1}N_{j(n)} = n^{-1}2^{j(n)} = o_n(1)$ and $e^{-v_n^2/4} \leq 1/j(n)$. \square

5.3. Proof of Proposition 3.1

We again restrict ourselves to the case of Gaussian errors ξ_i in (1.1). Recall that the vector $\widehat{\boldsymbol{\theta}}_j$ is defined as the subvector of $\widehat{\boldsymbol{\theta}}(j) = (\Psi(j)\Psi(j)^\top)^{-1} \Psi(j)\mathbf{Y}$, $j \leq j(n)$. The model equation (1.1) yields

$$\widehat{\boldsymbol{\theta}}(j) = \left(\Psi(j)\Psi(j)^\top \right)^{-1} \Psi(j)(f + \boldsymbol{\xi}) = \boldsymbol{\theta}(j) + \boldsymbol{\zeta}(j)$$

with $\boldsymbol{\theta}(j) = V(j)\Psi(j)f$ and $\boldsymbol{\zeta}(j) = V(j)\Psi(j)\boldsymbol{\xi}$ where $V(j) = (\Psi(j)\Psi(j)^\top)^{-1}$. Hence $\widehat{\boldsymbol{\theta}}_j = \boldsymbol{\theta}_j + \boldsymbol{\zeta}_j$ where $\boldsymbol{\theta}_j$ (resp. $\boldsymbol{\zeta}_j$) is the subvector of $\boldsymbol{\theta}(j)$ (resp. of $\boldsymbol{\zeta}(j)$) corresponding to the j th resolution level. This particularly implies that $\boldsymbol{\zeta}_j$ is a zero mean random vector with the covariance matrix V_j which is the submatrix of the matrix $V(j) = (\Psi(j)\Psi(j)^\top)^{-1}$. Moreover, if the errors ξ_i in (1.1) are Gaussian, then $\boldsymbol{\zeta}_j$ is for each $j \leq j(n)$ a Gaussian random vector with parameters $(0, V_j)$.

Let, for some $j \leq j(n)$, it holds

$$T_j^* = \frac{\|\boldsymbol{\theta}_j\|^2}{\sigma^2 \sqrt{2 \operatorname{tr} V_j^2}} \geq 3(\lambda_n^{1/2} + 1)^2 \quad (5.4)$$

with $\lambda_n = \max\{\lambda, 2\sqrt{\log j(n)}\}$. We shall show that under this condition it holds

$$\mathbf{P}_f(T_j < \lambda) \leq \delta(n) \rightarrow 0, \quad n \rightarrow \infty,$$

which obviously implies the assertion.

Observe first that

$$\begin{aligned} \mathbf{P}(T_j < \lambda) &= \mathbf{P}\left(\|\boldsymbol{\theta}_j + \boldsymbol{\zeta}_j\|^2 - \widehat{\sigma}^2 \operatorname{tr} V_j < \lambda \widehat{\sigma}^2 \sqrt{2 \operatorname{tr} V_j^2}\right) \\ &\leq \mathbf{P}\left(\|\boldsymbol{\theta}_j + \boldsymbol{\zeta}_j\|^2 - \sigma^2 \operatorname{tr} V_j < \lambda \sigma^2 \sqrt{2 \operatorname{tr} V_j^2} + (\widehat{\sigma}^2 - \sigma^2) \left(\lambda \sqrt{2 \operatorname{tr} V_j^2} + \operatorname{tr} V_j\right)\right) \\ &\leq \mathbf{P}\left(\|\boldsymbol{\theta}_j + \boldsymbol{\zeta}_j\|^2 - \sigma^2 \operatorname{tr} V_j - \|\boldsymbol{\theta}_j\|^2 < (\lambda + \lambda_n^{1/2}) \sigma^2 \sqrt{2 \operatorname{tr} V_j^2} - \|\boldsymbol{\theta}_j\|^2\right) \\ &\quad + \mathbf{P}\left((\widehat{\sigma}^2 - \sigma^2) \left(\lambda \sqrt{2 \operatorname{tr} V_j^2} + \operatorname{tr} V_j\right) < -\sigma^2 \lambda_n^{1/2} \sqrt{2 \operatorname{tr} V_j^2}\right). \end{aligned}$$

By Lemma 5.3 $\operatorname{tr} V_j (2 \operatorname{tr} V_j^2)^{-1/2} \leq \sqrt{N_j/2} \leq \sqrt{N_{j(n)}/2}$ for all $j \leq j(n)$. Further, by Lemma 5.2

$$\begin{aligned} \mathbf{P}\left(\frac{\widehat{\sigma}^2}{\sigma^2} - 1 < -\frac{\lambda_n^{1/2} \sqrt{2 \operatorname{tr} V_j^2}}{4\lambda \sqrt{2 \operatorname{tr} V_j^2} + \operatorname{tr} V_j}\right) \\ \leq 2 \exp\left(-\frac{\lambda_n n}{4(1+n^{-1})^2(\lambda + \sqrt{N_{j(n)}/2})^2}\right) = \delta_3(n) \end{aligned}$$

where $\delta_3(n) \rightarrow 0$ as $n \rightarrow \infty$ since $n/N_{j(n)} = n2^{-j(n)} \rightarrow \infty$.

Next, for every positive u , the inequality $\|\boldsymbol{\theta}\| \geq 3u$ implies $\|\boldsymbol{\theta}\|^2 - 2u\|\boldsymbol{\theta}\| - 3u^2 \geq 0$. Coupled with (5.4), this ensures, with $\tau_j = \sigma(2 \operatorname{tr} V_j^2)^{1/4}$ that

$$\begin{aligned} \|\boldsymbol{\theta}_j\|^2 &\geq \sqrt{4/3} \|\boldsymbol{\theta}_j\| (\lambda_n^{1/2} + 1) \tau_j + (\lambda_n^{1/2} + 1)^2 \tau_j^2 \\ &\geq \|\boldsymbol{\theta}_j\| (\lambda_n^{1/2} + 1) \tau_j + (\lambda_n + 2\lambda_n^{1/2} + 1) \tau_j^2. \end{aligned}$$

Now Proposition 6.2 with $\gamma = 1$ and $t = \lambda_n^{1/2}$ implies

$$\begin{aligned} \mathbf{P}(T_j < \lambda) &\leq \mathbf{P}\left(\|\boldsymbol{\theta}_j + \boldsymbol{\zeta}_j\|^2 - \sigma^2 \operatorname{tr} V_j - \|\boldsymbol{\theta}_j\|^2 < -(\lambda_n^{1/2} + 1) \|\boldsymbol{\theta}_j\| \tau_j - (\lambda_n^{1/2} + 1) \tau_j^2\right) + \delta_3(n) \\ &\leq 2e^{-\lambda_n^{1/2}/2} + \delta_3(n) \rightarrow 0, \quad n \rightarrow \infty \end{aligned}$$

as required.

5.4. Proof of Theorem 3.3

For the proof, we use the result of Proposition 3.1. Namely we show that the condition (3.2) of the theorem with \varkappa large enough contradict to the constraints

$$\|T_j^*\|^2 \leq t_n, \quad j \leq j(n), \quad (5.5)$$

with $t_n = 3 \left(1 + \max\{\lambda, 2\sqrt{\log j(n)}\}\right)^2$.

We begin by reduction of the problem of testing a linear hypothesis to the problem with a simple null hypothesis. Define coefficients θ_0, θ_1 by

$$(\theta_0, \theta_1) = \underset{(a,b)}{\operatorname{arginf}} \|f - a - b\psi_1\|_n = \underset{(a,b)}{\operatorname{arginf}} \sum_{i=1}^n (f(X_i) - a - bX_i)^2.$$

and set

$$f_0 = f - \theta_0 - \theta_1\psi_1.$$

Note that for all $j \geq 0$, the vectors $\boldsymbol{\theta}^*(j) = V(j)\Psi(j)f$ and $\boldsymbol{\theta}^*(j) = V(j)\Psi(j)f$ have the same components except the first two. Obviously the smoothness properties of f and f_0 also coincide and

$$\inf_{a,b} \|f - a - b\psi_1\|_n = \inf_{a,b} \|f_0 - a - b\psi_1\|_n$$

Recall also, that the linear trend in the regression function has no influence on our variance estimator $\hat{\sigma}^2$. Hence, replacing f by f_0 changes nothing in the test behaviour and we may suppose from the beginning that the coefficients θ_0^* and θ_1^* of the vector $\boldsymbol{\theta}^*(j)$ vanish.

About this new function f we know that

$$\begin{aligned} \|f\|_n &= \inf_{a,b} \|f - a - b\psi_1\|_n \geq \varrho(n), \\ &\inf_{g \in \mathcal{P}_s(j)} \|f - g\|_n = r_s(j), \end{aligned}$$

for all j from zero to $j(n)$.

Next we rewrite the constraints from (5.5) in term of the vectors $\|\boldsymbol{\theta}_j^*\|$, $j \leq j(n)$. Recall that $\boldsymbol{\theta}_j^*$ is the subvector of $\boldsymbol{\theta}^*(j)$ corresponding to j th level, and V_j is the corresponding submatrix of $V(j)$.

Let $\mathcal{L}(j)$ stand for the linear space generated by functions ψ_I , $I \in \mathcal{I}(j)$. We denote also by $\Pi(j)f$ the projection of f onto the space $\mathcal{L}(j)$ with respect to the norm $\|\cdot\|_n$,

$$\Pi(j)f = \underset{h \in \mathcal{L}(j)}{\operatorname{arginf}} \|f - h\|_n.$$

Particularly, $\Pi(0)f$ denotes the projection of f onto the space of linear functions (and hence, $\Pi(0)f = 0$) and, by definition of $\boldsymbol{\theta}(j)$,

$$\Pi(j)f = \sum_{I \in \mathcal{I}(j)} \theta_I^* \psi_I \tag{5.6}$$

where θ_I 's are the coefficients of the vector $\boldsymbol{\theta}^*(j)$.

Lemma 5.5. *For each $1 \leq j \leq j(n)$,*

$$\|\Pi(j)f\|_n \leq \|\Pi(j-1)f\|_n + \|\boldsymbol{\theta}_j^*\|.$$

Proof. Since $\mathcal{L}(j-1) \subseteq \mathcal{L}(j)$, then

$$\Pi(j-1)f = \Pi(j-1)\Pi(j)f.$$

When denoting $f(j) = \Pi(j)f$, one has $\Pi(j-1)f = \Pi(j-1)f(j)$ and we have to show that

$$\|\Pi(j-1)f(j)\|_n \geq \|f(j)\|_n - \|\theta_j^*\|.$$

In view of (5.6)

$$f(j) = \sum_{I \in \mathcal{I}(j)} \theta_I^* \psi_I.$$

Denote by f_j the part of this sum corresponding to the last level \mathcal{I}_j in $\mathcal{I}(j)$,

$$f_j = \sum_{I \in \mathcal{I}_j} \theta_I^* \psi_I.$$

By construction, the functions ψ_I , $I \in \mathcal{I}_j$, are orthonormal w.r.t. to the inner product $\|\cdot\|_n$ and particularly

$$\|f_j\|_n^2 = \sum_{I \in \mathcal{I}_j} |\theta_I^*|^2 = \|\theta_j^*\|^2.$$

Next, obviously $f(j) - f_j \in \mathcal{L}(j-1)$, and by definition of $\Pi(j)$,

$$\|f(j) - \Pi(j-1)f(j)\|_n \leq \|f(j) - (f(j) - f_j)\|_n = \|f_j\|_n = \|\theta_j^*\|$$

and the assertion follows by the triangle inequality. \square

Lemma 5.6. *Given $j \leq j(n)$, let (5.5) hold true for all $\ell \leq j$. Then*

$$\|\Pi(j)f\|_n^2 \leq \varkappa_1 C_V 2^{j/2} t_n$$

with $\varkappa_1 = 2^{1/2}(2^{1/4} - 1)^{-2}$.

Proof. Recursive application of Lemma 5.5 gives

$$\|\Pi(j)f\|_n \leq \sum_{\ell=0}^{j-1} \|\theta_\ell^*\|.$$

Here we have used that $\Pi(0)f = 0$. Now (5.5) and (D.iii) yield

$$\|\theta_\ell^*\|^2 \leq \sigma^2 t_n \sqrt{2 \operatorname{tr} V_\ell^2} \leq \sigma^2 t_n \sqrt{C_V^2 2^{\ell+1}}$$

and thus,

$$\|\Pi(j)f\|_n \leq \sum_{\ell=1}^j \left(2^{\ell/2} t_n C_V\right)^{1/2} = (C_V t_n)^{1/2} \sum_{\ell=1}^j 2^{\ell/4}$$

and the assertion follows by simple algebra. \square

Let now j_0 fulfill $2^{j_0} > s$ and $\mathcal{P}_s(j - j_0)$ denote the space of piecewise polynomials with piece length $2^{-(j-j_0)}$. Let now some $j \leq j(n)$ be fixed and let $g \in \mathcal{P}_s(j - j_0)$ be such that

$$\|f - g\|_n \leq r_s(j).$$

Lemma 5.7. *There is a constant $\varkappa_2 > 0$ depending on C_*, C^* and s only and such that for each j with $j_0 \leq j \leq j(n)$*

$$\|f\|_n \leq \varkappa_2 \{ \|\Pi(j)f\|_n + r_s(j) \}.$$

Proof. Let $g \in \mathcal{P}_s(j - j_0)$ be such that $\|f - g\|_n \leq r_s(j)$. Then

$$\|f\|_n \leq \|g\|_n + r_s(j)$$

and, since $\Pi(j)$ is a projector,

$$\begin{aligned} \|\Pi(j)f\|_n &= \|\Pi(j)g + \Pi(j)(f - g)\|_n \geq \|\Pi(j)g\|_n - \|\Pi(j)(f - g)\|_n \\ &\geq \|\Pi(j)g\|_n - r_s(j) \end{aligned}$$

and the assertion follows from

$$\|g\|_n^2 \leq \varkappa_3 \|\Pi(j)g\|_n^2.$$

Recall that g is a piecewise polynomial function on the partition A_I , $I \in \mathcal{I}_{j-j_0}$ and the projection $\Pi(j)g$ means the approximation of each polynomial on interval A_I of length $2^{-(j-j_0)}$ by piecewise constant functions with piece length 2^{-j} . Therefore, it suffices to prove that for each piece A_I and every polynomial $P(x) = a_0 + a_1x + \dots + a_{s-1}x^{s-1}$, it holds

$$\sum_{A_I} [\Pi(j)P(X_i)]^2 \geq \varkappa_3 \sum_{A_I} P^2(X_i)$$

where the constant \varkappa_3 depends on C_*, C^* and s only. The similar fact with integration instead of summation over the design points in A_I was stated in Ingster (1993) and we present here only a sketch of the proof for our situation.

The key idea of the proof can be formulated as a separate statement.

Lemma 5.8. *Let $P(x)$ be a polynomial of degree s and let m be an integer with $m > s + 1$. With $A_k = [(k-1)/m, k/m)$ for $k = 1, \dots, m$ Then for every measure μ on $[0, 1]$ with $0 < C_* \leq \mu(A_k) \leq C^* > 0$ for all $k \leq m$,*

$$\sum_{k=1}^m \left[\int_{A_k} P(x) \mu(dx) \right]^2 \geq \varkappa_3 \int_0^1 P^2(x) \mu(dx).$$

with a positive number \varkappa_3 depending on C_, C^* and s only.*

Proof. Let $a = (a_0, \dots, a_{s-1})$ be the vector of coefficients of P . Without loss of generality, we may assume that $\|a\|_\infty = \max_{j=0, \dots, s-1} \{|a_j|\} \leq 1$. Obviously, both

$$\begin{aligned} \|a\|_{\mu,1}^2 &= \left(\int_0^1 P(x) \mu(dx) \right)^2, \\ \|a\|_{\mu,2}^2 &= \int_0^1 P^2(x) \mu(dx) \end{aligned}$$

are scalar product in the space \mathbb{R}^s . Next, $\|a\|_{\mu,2} = 0$ only if $a = 0$ i.e. $P(x) \equiv 0$ and the same applies for $\|a\|_1$, since $P(x)$ has at most s roots and μ is supported on $m > s+1$ disjoint intervals. Note also that $\|a\|_{\mu,1}$ and $\|a\|_{\mu,2}$ are continuous functionals of a and μ and the space $\mathcal{M}_m(C_*, C^*)$ of measures μ on $[0, 1]$ satisfying the condition of the lemma is compact in the weak topology. Hence,

$$\sup_{a: \|a\|_\infty \leq 1} \sup_{\mu \in \mathcal{M}_m(C_*, C^*)} \frac{\|a\|_{\mu,2}}{\|a\|_{\mu,1}} = \varkappa_3 < \infty$$

as required. \square

Application of this result to each interval A_I , $I \in \mathcal{I}_{j-j_0}$ yields the desirable assertion. \square

The results of Lemma 5.5 through 5.7 yield the inequality

$$\|f\|_n \leq \varkappa_2 \left(r_s(j) + \sqrt{\varkappa_1 C_V 2^{j/2} \lambda_n} \right)$$

which contradicts to the constraints (5.5): $\|f\|_n \geq \varkappa \left(r_s(j) + \sqrt{2^{j/2} \lambda_n} \right)$ if \varkappa is large enough, and the theorem is proved.

5.5. Proof of Theorem 3.2

Now we disregard the assumption that the errors ξ_i in (1.1) are normally distributed and assume only they have 6 finite moments. We outline the proof of Theorem 3.2 only. Proposition 3.1 can be considered similarly.

Lemma 5.9. *Let the errors ξ_i in (1.1) are i.i.d. and satisfy $\mathbf{E}\xi_i = 0$, $\xi_i^2 = \sigma^2$ and $\mathbf{E}|\xi_i^2 - \sigma^2|^3 \leq C_6 \sigma^6$. Define $s_4^2 = 2\sigma^{-4} \mathbf{E}(\xi_1^2 - \sigma^2)^2$. Then, for $n \geq 36$ and every $\gamma > 0$,*

$$\mathbf{P} \left(\pm(\widehat{\sigma}^2 - \sigma^2) > (s_4 \gamma + \gamma + 1) \sqrt{\frac{4\sigma^4}{n}} \right) \leq 2e^{-\gamma^2/4} + rn^{-1/2}$$

where r depends on s_4 and C_6 only.

Proof. Similarly to the Gaussian case discussed in Section 5.2, it suffices to consider the case of the no-response model with the vanishing regression function. In this case, the

variance estimate $\hat{\sigma}^2$ is a quadratic form of the errors ξ_i which allows for the following representation:

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=2}^{n-1} \frac{(a_i \xi_{i-1} + b_i \xi_{i+1} - \xi_i)^2}{a_i^2 + b_i^2 + 1}$$

where $a_i = \frac{(X_{i+1} - X_i)}{(X_{i+1} - X_{i-1})}$, $b_i = \frac{(X_i - X_{i-1})}{(X_{i+1} - X_{i-1})}$, $i = 1, \dots, n$. The estimation error $\hat{\sigma}^2 - \sigma^2$ can be split into one centered diagonal quadratic form

$$\begin{aligned} Q_1 &= \frac{1}{n-2} \sum_{i=2}^{n-1} \frac{a_i^2 (\xi_{i-1}^2 - \sigma^2) + b_i^2 (\xi_{i+1}^2 - \sigma^2) + (\xi_i^2 - \sigma^2)}{a_i^2 + b_i^2 + 1} \\ &= \frac{1}{n-2} \sum_{i=2}^{n-1} \left(\frac{a_{i+1}^2}{a_{i+1}^2 + b_{i+1}^2 + 1} + \frac{1}{a_i^2 + b_i^2 + 1} + \frac{b_{i-1}^2}{a_{i-1}^2 + b_{i-1}^2 + 1} \right) (\xi_i^2 - \sigma^2) \end{aligned}$$

and one quadratic form Q_2 with vanishing diagonal elements:

$$\begin{aligned} Q_2 &= \frac{2}{n-2} \sum_{i=2}^{n-1} \frac{a_i b_i \xi_{i-1} \xi_{i+1} - a_i \xi_{i-1} \xi_i - b_i \xi_i \xi_{i+1}}{a_i^2 + b_i^2 + 1} \\ &= \frac{2}{n-2} \sum_{i=2}^{n-1} \frac{a_i b_i}{a_i^2 + b_i^2 + 1} \xi_{i-1} \xi_{i+1} \\ &\quad - \frac{2}{n-2} \sum_{i=2}^n \left(\frac{a_i}{a_i^2 + b_i^2 + 1} + \frac{b_{i-1}}{a_{i-1}^2 + b_{i-1}^2 + 1} \right) \xi_{i-1} \xi_i \end{aligned}$$

where $a_n = b_1 = 0$.

Obviously $Q_1 = \sum_{i=1}^n \lambda_i (\xi_i^2 - \sigma^2)$ with coefficients λ_i which, in view of the conditions $a_i, b_i \geq 0$, $a_i + b_i = 1$, fulfill $1/(2n-4) \leq \lambda_i \leq 2/(n-2)$. Since $\mathbf{E} (\xi_i^2 - \sigma^2)^2 \leq 2s_4^2 \sigma^4$, we bound

$$\mathbf{E} |Q_1|^2 = \sum_{i=1}^n \lambda_i^2 \mathbf{E} (\xi_i^2 - \sigma^2)^2 \leq 8s_4^2 n^{-1} \sigma^4.$$

Since also $\mathbf{E} |\xi_i^2 - \sigma^2|^3 \leq C_6 \sigma^6$, the Berry-Esseen theorem yields for every x , see Petrov (1975, Chapter 5)

$$\mathbf{P} \left(\frac{Q_1}{\sqrt{\mathbf{E} |Q_1|^2}} > x \right) - [1 - \Phi(x)] \leq r_1 C_6^{1/2} n^{-1/2}$$

where $\Phi(\cdot)$ is the Laplace function and r_1 is some absolute constant. This implies in view of $1 - \Phi(x) \leq e^{-x^2/2}$

$$\mathbf{P} \left(Q_1 > \gamma s_4 \sqrt{\frac{4\sigma^4}{n}} \right) \leq e^{-\gamma^2/4} + r_2 n^{-1/2}.$$

Further, similarly to Section 5.2, it holds $\mathbf{E} |Q_2|^2 \leq 4\sigma^4/n$. Moreover, it is easy to check that Q_2 fulfills the conditions of Proposition 6.3 with $G^2 = 4\sigma^4/n$ and some finite constant C_A and hence, by Corollary 6.1, with $\delta = 1$ and $\gamma > 0$,

$$\mathbf{P} (Q_2 > G(\gamma + 1)) \leq \mathbf{P} (\tilde{Q}_2 > G\gamma) + r_3 n^{-1/2}.$$

Finally, the quadratic form \tilde{Q}_2 of Gaussian random variables $\tilde{\xi}_i$ can be handled as in Section 5.2:

$$\mathbf{P} \left(\tilde{Q}_2 > \gamma \sqrt{\frac{4\sigma^4}{n}} \right) \leq e^{-\gamma^2/4}$$

if $n \geq 36$.

Combination of these results yields

$$\begin{aligned} & \mathbf{P} \left(\hat{\sigma}^2 - \sigma^2 > (s_4\gamma + \gamma + 1) \sqrt{\frac{4\sigma^4}{n}} \right) \\ & \leq \mathbf{P} \left(Q_1 > s_4\gamma \sqrt{\frac{4\sigma^4}{n}} \right) + \mathbf{P} \left(Q_2 > (\gamma + 1) \sqrt{\frac{4\sigma^4}{n}} \right) \\ & \leq 2e^{-\gamma^2/4} + rn^{-1/2}. \end{aligned}$$

Similarly one can get an upper bound for $-(\hat{\sigma}^2 - \sigma^2)$.

□

In the same way one can extend the result of Lemma 5.2 to the non-Gaussian case: $\hat{\sigma}^2$ estimates the true variance σ^2 at the rate $n^{-1/2}$ provided that f is sufficiently smooth.

Now we turn to Theorem 3.2. It obviously suffices to show that the distribution of the test statistic T^* can be approximated by a similar distribution corresponding to the case of Gaussian errors. Then the result follows from Theorem 3.1.

As in the proof of Theorem 3.1, the general case can be reduced to the no-response model with the vanishing regression function. Further, since the difference $\hat{\sigma}^2 - \sigma^2$ is of order $n^{-1/2}$, it suffices to consider the expressions T'_j , $j \leq j(n)$, defined by

$$T'_j = \frac{1}{\sqrt{2\sigma^4 \operatorname{tr} V_j^2}} \left(\sum_{I \in \mathcal{I}_j} |\hat{\theta}_I|^2 - \sigma^2 \operatorname{tr} V_j \right) = \frac{S_j - \sigma^2 \operatorname{tr} V_j}{\sqrt{2\sigma^4 \operatorname{tr} V_j^2}}$$

where $\hat{\theta}_I$ are elements of the vector $\boldsymbol{\theta}(j)$, cf. the proof of Lemma 5.4. Under the no-response hypothesis, this vector admits the representation, $\boldsymbol{\theta}(j) = W(j)\boldsymbol{\xi}$ with $W(j) = (\Psi(j)^\top \Psi(j))^{-1} \Psi(j)^\top$, see (2.6). If P_j denotes the mapping from $\mathcal{I}(j)$ into \mathcal{I}_j , then $\boldsymbol{\theta}_j = P_j \boldsymbol{\theta}(j) = P_j W(j)\boldsymbol{\xi}$ and

$$S_j = \|\hat{\boldsymbol{\theta}}_j\|^2 = \boldsymbol{\xi}^\top W(j)^\top P_j^\top P_j W(j)\boldsymbol{\xi} = \boldsymbol{\xi}^\top A_j \boldsymbol{\xi}_j$$

with $A_j = W(j)^\top P_j^\top P_j W(j)$, so that S_j is a quadratic form of the errors ξ_i . We also know that $V_j = P_j W(j) W(j)^\top P_j^\top$, and $\mathbf{E} S_j = \sigma^2 \operatorname{tr} A_j = \sigma^2 \operatorname{tr} V_j$. Moreover, see (2.8), under Gaussian errors ξ_i , it also holds $\mathbf{E} (S_j - \mathbf{E} S_j)^2 = \sigma^4 \operatorname{tr} V_j$. Hence, each of T'_j is a centered and normalized quadratic form of ξ_i 's. This form in turns can be represented as a sum of a diagonal form $T_j^{(1)}$ and a quadratic form $T_j^{(2)}$ with vanishing diagonal

terms. We first show that the impact of diagonal terms is negligible and then apply Corollary 6.2 to $T_j^{(2)}$'s.

Let o_i denote the i -th basis vector in \mathbb{R}^n . Then the i -th diagonal element a_{ii} of A_j is equal to $o_i^\top A_j o_i$:

$$\begin{aligned} a_{ii} &= o_i^\top A_j o_i \\ &= o_i^\top \Psi(j)^\top \left(\Psi(j)^\top \Psi(j) \right)^{-1} P_j^\top P_j \left(\Psi(j)^\top \Psi(j) \right)^{-1} P_j^\top P_j \left(\Psi(j)^\top \Psi(j) \right)^{-1} \Psi(j) o_i. \end{aligned}$$

Clearly

$$\left\| \left(\Psi(j)^\top \Psi(j) \right)^{-1} P_j^\top P_j \left(\Psi(j)^\top \Psi(j) \right)^{-1} \right\| \leq \left\| \left(\Psi(j)^\top \Psi(j) \right)^{-2} \right\| = \|V(j)^2\| \leq C_V^2.$$

Next, for every Haar level $\ell \leq j$, there exists only one index $I \in \mathcal{I}_\ell$ such that $\psi_I(X_i) \neq 0$. More precisely, for this index I , it holds $\psi_I(X_i) = \pm 1/\sqrt{M_I}$ where M_I is the number of design points in the interval A_I corresponding to the index I . Condition (D.i) implies $M_I \geq C_* n 2^{-\ell}$ for every $I \in \mathcal{I}_\ell$. Also $\psi_0(X_i) = n^{-1/2}$ and $\psi_1(X_i) = X_i \left(\sum_{i'=1}^n X_{i'}^2 \right)^{-1/2}$. Hence, the definition of the matrix $\Psi(j)$ and condition (D.i) provide

$$|\Psi(j) o_i| \leq n^{-1/2} + \left(\sum_{i'=1}^n X_{i'}^2 \right)^{-1/2} + \sum_{\ell=0}^j \sqrt{\frac{2^\ell}{nC_*}} < 3C_*^{-1/2} 2^{j/2} n^{-1/2}. \quad (5.7)$$

Therefore,

$$\begin{aligned} a_{ii} &\leq |\Psi(j) o_i|^2 \left\| \left(\Psi(j)^\top \Psi(j) \right)^{-1} P_j^\top P_j \left(\Psi(j)^\top \Psi(j) \right)^{-1} \right\| \\ &\leq 9C_*^{-1} 2^j n^{-1} C_V^2. \end{aligned}$$

Define $G_j^2 = \sigma^4 \text{tr} A_j^2$. Note

$$\begin{aligned} \text{tr} A_j^2 &= \text{tr} W(j)^\top P_j^\top P_j W(j) W(j)^\top P_j^\top P_j W(j) \\ &= \text{tr} P_j W(j) W(j)^\top P_j^\top P_j W(j) W(j)^\top P_j^\top \\ &= \text{tr} V_j^2 \end{aligned}$$

so that $T_j^{(1)} = G_j^{-1} \sum_{i=1}^n a_{ii}(\xi_i^2 - \sigma^2)$. The condition (D.ii) implies $\text{tr} A_j^2 \geq C_D 2^j$. Now, for every $\delta > 0$,

$$\begin{aligned}
\mathbf{P} \left(\max_{j=0, \dots, j(n)} T_j^{(1)} > \delta \right) &\leq \sum_{j=0}^{j(n)} \mathbf{P} \left(T_j^{(1)} > \delta \right) \\
&\leq \delta^{-2} \sum_{j=0}^{j(n)} \mathbf{E} \left| T_j^{(1)} \right|^2 \\
&\leq \delta^{-2} \sum_{j=0}^{j(n)} 2G_j^{-2} \sigma^4 a_{ii}^2 \\
&\leq 2\delta^{-2} \sum_{j=0}^{j(n)} 2C_D^{-1} 2^{-j} n (9C_*^{-1} 2^j n^{-1} C_V^2)^2 \\
&\leq C\delta^{-2} n^{-1} 2^{j(n)+1} \rightarrow 0, \quad n \rightarrow \infty.
\end{aligned}$$

Next we consider $T_j^{(2)}$ which is obtained from T_j' by removing the diagonal terms. This quadratic form can be approximated (in distribution) by a similar one with Gaussian errors $\tilde{\xi}_i$ at a reasonable rate provided that the corresponding value C_A , defined as n times the ratio of the maximal diagonal element of the matrix $\sigma^4 A_j^2$ to $G_j^2 = \sigma^4 \text{tr} A_j^2$, see (6.2) and Remark 6.1, remains bounded.

The i -th diagonal element d_i of A_j^2 is equal to $o_i^\top A_j^2 o_i$:

$$\begin{aligned}
d_i &= o_i^\top A_j^2 o_i \\
&= o_i^\top \left\{ \Psi(j) \left(\Psi(j)^\top \Psi(j) \right)^{-1} P_j^\top P_j \left(\Psi(j)^\top \Psi(j) \right)^{-1} \Psi(j)^\top \right\}^2 o_i \\
&= o_i^\top \Psi(j)^\top \left(\Psi(j)^\top \Psi(j) \right)^{-1} P_j^\top P_j \left(\Psi(j)^\top \Psi(j) \right)^{-1} P_j^\top P_j \left(\Psi(j)^\top \Psi(j) \right)^{-1} \Psi(j) o_i.
\end{aligned}$$

Clearly

$$\begin{aligned}
&\left\| \left(\Psi(j)^\top \Psi(j) \right)^{-1} P_j^\top P_j \left(\Psi(j)^\top \Psi(j) \right)^{-1} P_j^\top P_j \left(\Psi(j)^\top \Psi(j) \right)^{-1} \right\| \\
&\leq \left\| \left(\Psi(j)^\top \Psi(j) \right)^{-3} \right\| = \|V(j)^3\| \leq C_V^3.
\end{aligned}$$

The use of (5.7) provides

$$\begin{aligned}
d_i &\leq |\Psi(j) o_i|^2 \left\| \left(\Psi(j)^\top \Psi(j) \right)^{-1} P_j^\top P_j \left(\Psi(j)^\top \Psi(j) \right)^{-1} P_j^\top P_j \left(\Psi(j)^\top \Psi(j) \right)^{-1} \right\| \\
&\leq 9C_*^{-1} 2^j n^{-1} C_V^3
\end{aligned}$$

and

$$C_A \leq \frac{9C_*^{-1} C_V^3 2^j}{C_D 2^j} = \frac{9C_*^{-1} C_V^3}{C_D}$$

that is, the value C_A is bounded by a fixed constant depending on design regularity only.

By Corollary 6.2, the joint distribution of $T_j^{(2)}$, $j \leq j(n)$, and the distribution of their maximum, can be approximated by the distribution of similar quadratic forms of Gaussian r.v.'s which implies the required assertion.

6. Appendix

Here we discuss briefly some general properties of quadratic forms of random variables. We first consider the case when the underlying random variables are Gaussian and establish an exponential bound for deviations of such forms over certain level. Next we show how an arbitrary quadratic form of independent random variables can be approximated (in distribution) by a similar quadratic form of Gaussian random variables.

6.1. Deviation probabilities for quadratic forms of Gaussian random variables

Let $\varepsilon_1, \dots, \varepsilon_N$ be Gaussian random variables with zero mean and the covariance $N \times N$ matrix V , i.e. $V = \mathbf{E}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top$ where $\boldsymbol{\varepsilon}$ denotes the vector $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)^\top$.

We first present the following general results about quadratic forms of Gaussian random variables.

Proposition 6.1. *Let $\varepsilon_1, \dots, \varepsilon_N$ be Gaussian random variables with zero mean and the covariance matrix $V := \mathbf{E}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top$. Then*

$$\begin{aligned} \mathbf{E}\|\boldsymbol{\varepsilon}\|^2 &:= \mathbf{E}(\varepsilon_1^2 + \dots + \varepsilon_N^2) = \text{tr } V, \\ \mathbf{E}(\|\boldsymbol{\varepsilon}\|^2 - \text{tr } V)^2 &= 2 \text{tr } V^2. \end{aligned}$$

Moreover, for $\gamma \leq \|V\|^{-1}\sqrt{\text{tr } V^2/2}$ and each $t \geq 0$,

$$\mathbf{P}\left(\pm(\|\boldsymbol{\varepsilon}\|^2 - \text{tr } V) > (\gamma + t)\sqrt{2 \text{tr } V^2}\right) \leq e^{-\gamma t/2 - \gamma^2/4},$$

Proof. Let $V = U^\top \Lambda U$ be a diagonal representation of V with a diagonal matrix $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_N\}$ and an ortonormal matrix U . It is well known that $\zeta = \Lambda^{-1/2}U\boldsymbol{\varepsilon}$ is a standard Gaussian vector and $\|\boldsymbol{\varepsilon}\|^2 = \zeta^\top \Lambda \zeta$. Also it holds $\text{tr } V = \lambda_1 + \dots + \lambda_N$, $\text{tr } V^2 = \lambda_1^2 + \dots + \lambda_N^2$ and $\|V\| = \max\{\lambda_1, \dots, \lambda_N\}$. To bound the expression $\|\boldsymbol{\varepsilon}\|^2 - \text{tr } V$, we apply the exponential Chebyshev inequality: with each $\mu \geq 0$ satisfying $2\mu\lambda_i < 1$ and

every z

$$\begin{aligned}
\mathbf{P}(\|\boldsymbol{\varepsilon}\|^2 - \text{tr } V > z) &\leq e^{-\mu z} \mathbf{E} \exp\{\mu(\|\boldsymbol{\varepsilon}\|^2 - \text{tr } V)\} \\
&= e^{-\mu z} \mathbf{E} \exp\left\{\mu \sum_{i=1}^N \lambda_i (\zeta_i^2 - 1)\right\} \\
&= e^{-\mu z} \prod_{i=1}^N \mathbf{E} \exp\{\mu \lambda_i (\zeta_i^2 - 1)\} \\
&= \exp\left\{-\mu z - \mu \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \frac{1}{2} \log(1 - 2\mu \lambda_i)\right\}.
\end{aligned}$$

We now set $\mu = \frac{z}{2\sqrt{2 \text{tr } V^2}}$ so that $2\mu \lambda_i = \frac{z \lambda_i}{\sqrt{2 \text{tr } V^2}} < 1/2$ and use that $-\log(1-u) \leq u+u^2$ for $0 \leq u \leq 1/2$. This yields

$$\begin{aligned}
\mathbf{P}\left(\|\boldsymbol{\varepsilon}\|^2 - \text{tr } V > (\gamma + t)\sqrt{2 \text{tr } V^2}\right) &\leq \exp\left(-\frac{\gamma(\gamma + t)}{2} + \frac{\gamma^2}{4 \text{tr } V^2} \sum_{i=1}^N \lambda_i^2\right) \\
&= \exp(-\gamma t/2 - \gamma^2/4)
\end{aligned}$$

as required. The bound for $-(\|\boldsymbol{\varepsilon}\|^2 - \text{tr } V)$ is proved in the same line. \square

Further, for a deterministic vector $\mathbf{c} = (c_1, \dots, c_N)^\top$ from \mathbb{R}^N , we consider quadratic forms of type

$$\|\mathbf{c} + \boldsymbol{\varepsilon}\|^2 = \sum_{j=1}^N |c_j + \varepsilon_j|^2.$$

Proposition 6.2. *Let $\varepsilon_1, \dots, \varepsilon_N$ be Gaussian random variables with zero mean and the covariance matrix $\sigma^2 V$. Then it holds for any vector $\mathbf{c} = (c_1, \dots, c_N)^\top$ in \mathbb{R}^N*

$$\begin{aligned}
\mathbf{E}\|\mathbf{c} + \boldsymbol{\varepsilon}\|^2 &= \|\mathbf{c}\|^2 + \text{tr } V, \\
\mathbf{E}\left(\|\mathbf{c} + \boldsymbol{\varepsilon}\|^2 - \|\mathbf{c}\|^2 - \text{tr } V\right)^2 &= 4\mathbf{c}^\top V \mathbf{c} + 2 \text{tr } V^2,
\end{aligned}$$

Moreover, for every positive γ with $\gamma \leq \|V\|^{-1} \sqrt{\text{tr } V^2/2}$ and every $t \geq 0$

$$\mathbf{P}\left(\pm(\|\mathbf{c} + \boldsymbol{\varepsilon}\|^2 - \|\mathbf{c}\|^2 - \text{tr } V) > \gamma \|\mathbf{c}\| (2 \text{tr } V^2)^{1/4} + (\gamma + t)\sqrt{2 \text{tr } V^2}\right) \leq 2e^{-\gamma^2/4 - \gamma t/2}.$$

Proof. With vector notation, the studied quadratic form can be rewritten as $\|\mathbf{c} + \boldsymbol{\varepsilon}\|^2 = (\mathbf{c} + \boldsymbol{\varepsilon})^\top (\mathbf{c} + \boldsymbol{\varepsilon})$. Now, since $\mathbf{E}\boldsymbol{\varepsilon}_i = 0$, it holds

$$\mathbf{E}\|\mathbf{c} + \boldsymbol{\varepsilon}\|^2 = \mathbf{E}\left(\|\mathbf{c}\|^2 + 2\mathbf{c}^\top \boldsymbol{\varepsilon} + \|\boldsymbol{\varepsilon}\|^2\right) = \|\mathbf{c}\|^2 + \mathbf{E}\|\boldsymbol{\varepsilon}\|^2 = \|\mathbf{c}\|^2 + \text{tr } V.$$

Next,

$$\begin{aligned}
\text{Var } \|\mathbf{c} + \boldsymbol{\varepsilon}\|^2 &= \mathbf{E}\left(\|\mathbf{c} + \boldsymbol{\varepsilon}\|^2 - \mathbf{E}\|\mathbf{c} + \boldsymbol{\varepsilon}\|^2\right)^2 \\
&= \mathbf{E}\left(2\mathbf{c}^\top \boldsymbol{\varepsilon} + \|\boldsymbol{\varepsilon}\|^2 - \text{tr } V\right)^2 \\
&= 4\mathbf{E}|\mathbf{c}^\top \boldsymbol{\varepsilon}|^2 + 4\mathbf{E}\mathbf{c}^\top \boldsymbol{\varepsilon} (\|\boldsymbol{\varepsilon}\|^2 - \text{tr } V) + \mathbf{E}(\|\boldsymbol{\varepsilon}\|^2 - \text{tr } V)^2.
\end{aligned}$$

The Gaussian vector $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, V)$ fulfills

$$\begin{aligned} \mathbf{E} \boldsymbol{\varepsilon} (\|\boldsymbol{\varepsilon}\|^2 - \text{tr } V) &= 0, \\ \mathbf{E} |\mathbf{c}^\top \boldsymbol{\varepsilon}|^2 &= \mathbf{c}^\top (\mathbf{E} \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top) \mathbf{c} = \mathbf{c}^\top V \mathbf{c} \end{aligned}$$

so that in view of Lemma 6.1 $\text{Var} \|\mathbf{c} + \boldsymbol{\varepsilon}\|^2 = 4\mathbf{c}^\top V \mathbf{c} + 2 \text{tr } V^2$ as required.

Let now $\gamma \geq 1$ be fixed such that $\gamma \leq \|V\|^{-1} \sqrt{\text{tr } V^2/2}$. This particularly means that $\|V\| \leq \sqrt{\text{tr } V^2/2}$. Note that the scalar product $\mathbf{c}^\top \boldsymbol{\varepsilon}$ is a linear combination of the Gaussian zero mean random variables and it is therefore Gaussian as well with $\mathbf{E} \mathbf{c}^\top \boldsymbol{\varepsilon} = 0$ and $\mathbf{E} |\mathbf{c}^\top \boldsymbol{\varepsilon}|^2 = \mathbf{c}^\top V \mathbf{c}$. This yields for every $\gamma \geq 1$

$$\mathbf{P} \left(\mathbf{c}^\top \boldsymbol{\varepsilon} > \gamma \sqrt{\mathbf{c}^\top V \mathbf{c}} \right) \leq e^{-\gamma^2/2}.$$

The condition $\|V\| \leq \sqrt{\text{tr } V^2/2}$ provides $\mathbf{c}^\top V \mathbf{c} \leq \|\mathbf{c}\|^2 \|V\| \leq \|\mathbf{c}\|^2 \sqrt{\text{tr } V^2/2}$. Combining this inequality with the previous one implies

$$\mathbf{P} \left(2\mathbf{c}^\top \boldsymbol{\varepsilon} > (\gamma + t) \|\mathbf{c}\| (2 \text{tr } V^2)^{1/4} \right) \leq e^{-(\gamma+t)^2/4}.$$

Next, by Lemma 6.1

$$\mathbf{P} \left(\|\boldsymbol{\varepsilon}\|^2 - \text{tr } V > (\gamma + t) \sqrt{2 \text{tr } V^2} \right) \leq e^{-\gamma^2/4 - \gamma t/2}.$$

Summing up the previous estimates, we obtain

$$\begin{aligned} &\mathbf{P} \left(\sum_{j=1}^N |c_j + \varepsilon_j|^2 - \text{tr } V > \|\mathbf{c}\|^2 + (\gamma + t) \|\mathbf{c}\| (2 \text{tr } V^2)^{1/4} + (\gamma + t) \sqrt{2 \text{tr } V^2} \right) \\ &= \mathbf{P} \left(2\mathbf{c}^\top \boldsymbol{\varepsilon} + \|\boldsymbol{\varepsilon}\|^2 - \text{tr } V > (\gamma + t) \|\mathbf{c}\| (2 \text{tr } V^2)^{1/4} + (\gamma + t) \sqrt{2 \text{tr } V^2} \right) \\ &\leq \mathbf{P} \left(2\mathbf{c}^\top \boldsymbol{\varepsilon} > (\gamma + t) \|\mathbf{c}\| (2 \text{tr } V^2)^{1/4} \right) + \mathbf{P} \left(\|\boldsymbol{\varepsilon}\|^2 - \text{tr } V > (\gamma + t) \sqrt{2 \text{tr } V^2} \right) \\ &\leq 2e^{-\gamma^2/4 - \gamma t/2} \end{aligned}$$

as required. \square

6.2. Gaussian approximation for quadratic forms

In what follows we consider quadratic forms $\sum_{i=1}^n \sum_{\ell=1}^n a_{i\ell} \xi_i \xi_\ell$ of independent but not necessarily normal random variables ξ_1, \dots, ξ_n with vanishing diagonal coefficients, i.e. $a_{ii} = 0$. We aim to show that, under moment conditions on ξ_i 's and mild assumptions on the coefficients of the quadratic form, the asymptotic distribution of this quadratic form only weakly depends on the particular distribution of ξ_i 's and, as a consequence, it can be approximated by a distribution of a similar quadratic form of Gaussian r.v.'s with the same first and second moments.

Let $A = (a_{i\ell}, i, j = 1, \dots, n)$ be a $n \times n$ symmetric matrix with $a_{ii} = 0$ for all i , and let ξ_1, \dots, ξ_n be independent zero mean r.v.'s with $\mathbf{E} \xi_i^4 < \infty$ for all i . Define $\sigma_i^2 = \mathbf{E} \xi_i^2$. We study some properties of the quadratic form $\sum_{i=1}^n \sum_{j=1}^n a_{ij} \xi_i \xi_j$.

Lemma 6.1. *It holds*

$$\begin{aligned} \mathbf{E} \sum_{i=1}^n \sum_{\ell=1}^n a_{i\ell} \xi_i \xi_\ell &= \sum_{i=1}^n a_{ii} \sigma_i^2 = 0, \\ \mathbf{E} \left\{ \sum_{i=1}^n \sum_{\ell=1}^n a_{i\ell} \xi_i \xi_\ell \right\}^2 &= 2 \sum_{i=1}^n \sum_{\ell \neq i}^n a_{i\ell}^2 \sigma_i^2 \sigma_\ell^2. \end{aligned} \quad (6.1)$$

Proof. Obvious. Here it is only important that the diagonal elements a_{ii} vanish. \square

By $A(\xi_1, \dots, \xi_n)$ we denote the corresponding quadratic form, that is

$$A(\xi_1, \dots, \xi_n) = \sum_{i=1}^n \sum_{\ell \neq i}^n a_{i\ell} \xi_i \xi_\ell.$$

Let also $\tilde{\xi}_1, \dots, \tilde{\xi}_n$ be a sequence of independent Gaussian r.v.'s with $\mathbf{E}\tilde{\xi}_i = 0$ and $\mathbf{E}\tilde{\xi}_i^2 = \sigma_i^2$, $i = 1, \dots, n$. Define another quadratic form

$$A(\tilde{\xi}_1, \dots, \tilde{\xi}_n) = \sum_{i=1}^n \sum_{\ell \neq i}^n a_{i\ell} \tilde{\xi}_i \tilde{\xi}_\ell$$

Clearly $\mathbf{E}A(\tilde{\xi}_1, \dots, \tilde{\xi}_n) = 0$ and $\mathbf{E}|A(\tilde{\xi}_1, \dots, \tilde{\xi}_n)|^2 = \mathbf{E}|A(\xi_1, \dots, \xi_n)|^2$.

Proposition 6.3. *Let $\mathbf{E}\xi_i^4 \leq C_4\sigma_i^4$ for some fixed constant $C_4 \geq 3$. Let, for a symmetric matrix A with $a_{ii} = 0$ for $i = 1, \dots, n$, and for a normalizing constant G , the numbers C_A be defined by*

$$C_A = \max_{i=1, \dots, n} nG^{-2} \sum_{\ell=1}^n a_{i\ell}^2 \sigma_i^2 \sigma_\ell^2. \quad (6.2)$$

Then, for every three times continuously differentiable function f , it holds

$$\left| \mathbf{E}f(G^{-1}A(\xi_1, \dots, \xi_n)) - \mathbf{E}f(G^{-1}A(\tilde{\xi}_1, \dots, \tilde{\xi}_n)) \right| \leq \frac{8}{3} f_3(C_A C_A)^{3/2} n^{-1/2}$$

where f_3 means the maximum of the absolute value of the third derivative of f , that is, $f_3 = \sup_x |f'''(x)|$.

Remark 6.1. The value C_A can be easily evaluated for the case of an homogeneous noise when all σ_i^2 coincide with some σ^2 . Clearly each sum $d_i = \sum_{\ell=1}^n a_{i\ell}^2$ is i -th diagonal element of A^2 and $C_A \leq G^{-2} \max_{i=1, \dots, n} \{nd_i\}$.

Remark 6.2. The conditions of Proposition 6.3 do not guarantee that the distribution of $G^{-1}A(\xi_1, \dots, \xi_n)$ is close to some normal distribution. A typical example which just meets in hypothesis testing framework corresponds to the quadratic form $A(\xi_1, \dots, \xi_n) = (\xi_1 + \dots + \xi_n)^2$. which, even with normal ξ_i 's, has the χ^2 -distribution.

Proof. The change ξ_i for ξ_i/σ_i and $a_{i\ell}$ for $a_{i\ell}\sigma_i\sigma_\ell$ allows to reduce the general case to the situation with $\sigma_i = 1$ for all i . Hence, for the sake of notation simplicity, we suppose that $\xi_i^2 = 1$, $i = 1, \dots, n$.

We use the following obvious inequality

$$\begin{aligned} & \left| \mathbf{E}f\left(G^{-1}A(\xi_1, \dots, \xi_n)\right) - \mathbf{E}f\left(G^{-1}A(\tilde{\xi}_1, \dots, \tilde{\xi}_n)\right) \right| \\ & \leq \sum_{i=1}^n \left| \mathbf{E}f\left(G^{-1}A(\xi_1, \dots, \xi_i, \tilde{\xi}_{i+1}, \dots, \tilde{\xi}_n)\right) - \mathbf{E}f\left(G^{-1}A(\xi_1, \dots, \xi_{i-1}, \tilde{\xi}_i, \dots, \tilde{\xi}_n)\right) \right| \end{aligned} \quad (6.3)$$

where we assume $\xi_0 = \tilde{\xi}_{n+1} = 0$. We evaluate the last summand here, all others can be bounded in the same way. Denote

$$\begin{aligned} u_{n-1} &= G^{-1} \sum_{i=1}^{n-1} \sum_{\ell \neq i}^{n-1} a_{i\ell} \xi_i \xi_\ell, \\ \Delta_n &= G^{-1}A(\xi_1, \dots, \xi_n) - u_{n-1} = 2G^{-1}\xi_n \sum_{i=1}^{n-1} a_{in} \xi_i, \\ \tilde{\Delta}_n &= G^{-1}A(\xi_1, \dots, \xi_{n-1}, \tilde{\xi}_n) - u_{n-1} = 2G^{-1}\tilde{\xi}_n \sum_{i=1}^{n-1} a_{in} \xi_i. \end{aligned}$$

The Taylor expansion yields

$$\begin{aligned} & \left| \mathbf{E}f\left(G^{-1}A(\xi_1, \dots, \xi_n)\right) - \mathbf{E}f\left(G^{-1}A(\xi_1, \dots, \xi_{n-1}, \tilde{\xi}_n)\right) \right| \\ & \leq \left| \mathbf{E}f'(u_{n-1})(\Delta_n - \tilde{\Delta}_n) \right| + \frac{1}{2} \left| \mathbf{E}f''(u_{n-1})(\Delta_n^2 - \tilde{\Delta}_n^2) \right| + \frac{f_3}{6} (\mathbf{E}|\Delta_n|^3 + \mathbf{E}|\tilde{\Delta}_n|^3). \end{aligned} \quad (6.4)$$

Since ξ_n and $\tilde{\xi}_n$ are independent of ξ_1, \dots, ξ_{n-1} and since $\mathbf{E}\xi_n = \mathbf{E}\tilde{\xi}_n = 0$, $\mathbf{E}\xi_n^2 = \mathbf{E}\tilde{\xi}_n^2 = 1$, taking the conditional expectation given ξ_1, \dots, ξ_{n-1} , we obtain

$$\mathbf{E}\left(\Delta_n - \tilde{\Delta}_n \mid \xi_1, \dots, \xi_{n-1}\right) = 0 \quad (6.5)$$

$$\mathbf{E}\left(\Delta_n^2 - \tilde{\Delta}_n^2 \mid \xi_1, \dots, \xi_{n-1}\right) = 0. \quad (6.6)$$

Further we evaluate $\mathbf{E}|\Delta_n|^3$ and $\mathbf{E}|\tilde{\Delta}_n|^3$. Note first that, since $\mathbf{E}\xi_n^4 \leq C_4$ with $C_4 \geq 3$,

$$\begin{aligned} \mathbf{E}\left(\sum_{i=1}^{n-1} a_{in} \xi_i\right)^4 &= \sum_{i=1}^{n-1} a_{in}^4 \mathbf{E}\xi_i^4 + 3 \sum_{\ell \neq i}^{n-1} a_{in}^2 a_{\ell n}^2 \\ &\leq \sum_{i=1}^{n-1} a_{in}^4 (C_4 - 3) + 3 \left(\sum_{i=1}^{n-1} a_{in}^2\right)^2 \\ &\leq C_4 \left(\sum_{i=1}^{n-1} a_{in}^2\right)^2. \end{aligned}$$

Now the Hölder inequality yields in view of $\mathbf{E}|\xi_n|^3 \leq C_4^{3/4}$

$$\begin{aligned} G^3 \mathbf{E}|\Delta_n|^3 &= \mathbf{E}|\xi_n|^3 \mathbf{E} \left| 2 \sum_{i=1}^{n-1} a_{in} \xi_i \right|^3 \\ &\leq 8C_4^{3/4} \left\{ \mathbf{E} \left(\sum_{i=1}^{n-1} a_{in} \xi_i \right)^4 \right\}^{3/4} \\ &\leq 8C_4^{3/2} \left(\sum_{i=1}^n a_{in}^2 \right)^{3/2} \end{aligned}$$

and the condition $G^{-2} \sum_{i=1}^n a_{in}^2 \leq n^{-1} C_A$ provides

$$\mathbf{E}|\Delta_n|^3 \leq 8(C_4 C_A)^{3/2} n^{-3/2}. \quad (6.7)$$

For the Gaussian r.v. $s_n \tilde{\xi}_n$, the similar bound applies:

$$\mathbf{E}|\tilde{\Delta}_n|^3 \leq 8(C_4 C_A)^{3/2} n^{-3/2}. \quad (6.8)$$

Substituting these estimates as well as (6.5) and (6.6) in (6.4) implies

$$\left| \mathbf{E}f \left(\frac{A(\xi_1, \dots, \xi_n)}{G} \right) - \mathbf{E}f \left(\frac{A(\xi_1, \dots, \xi_{n-1}, \tilde{\xi}_n)}{G} \right) \right| \leq \frac{16}{6} f_3(C_4 C_A)^{3/2} n^{-3/2}.$$

Similar bounds hold for the other summands in (6.3). Summing them out, we obtain

$$\left| \mathbf{E}f(G^{-1}A(\xi_1, \dots, \xi_n)) - \mathbf{E}f(G^{-1}A(\tilde{\xi}_1, \dots, \tilde{\xi}_n)) \right| \leq \frac{8}{3} f_3(C_4 C_A)^{3/2} n^{-1/2}$$

as required. \square

Corollary 6.1. *Under the conditions of Proposition 6.3, for each $\delta > 0$ and every x*

$$\mathbf{P}(G^{-1}A(\xi_1, \dots, \xi_n) > x) \leq \mathbf{P}(G^{-1}A(\tilde{\xi}_1, \dots, \tilde{\xi}_n) > x - \delta) + \text{Const.} C_A^{3/2} n^{-1/2} \delta^{-3}$$

with a constant Const. depending on C_4 only. If, in addition, $G^2 \geq \mathbf{E}|A(\xi_1, \dots, \xi_n)|^2$, then

$$\mathbf{P}(G^{-1}A(\xi_1, \dots, \xi_n) > x) \leq \mathbf{P}(G^{-1}A(\tilde{\xi}_1, \dots, \tilde{\xi}_n) > x) + \text{Const.} C_A^{3/2} n^{-1/2} \delta^{-3} + \delta.$$

Proof. Let a smooth function f fulfill $f(u) = 0$ for $u \leq -1$ and $f(u) = 1$ for $u \geq 0$. Define $C_f = \sup_u |f'''(u)|$. Now, given x and $\delta > 0$, set $f_{x,\delta}(u) = f(\delta^{-1}(u - x))$. Obviously $f_{x,\delta}(u) = 0$ for $u \leq x - \delta$ and $f_{x,\delta}(u) = 1$ for $u \geq x$ and also $|f_{x,\delta}'''(u)| \leq C_f \delta^{-3}$.

Next, by Proposition 6.3

$$\begin{aligned} \mathbf{P}(G^{-1}A(\xi_1, \dots, \xi_n) > x) &\leq \mathbf{E}f_{x,\delta}(G^{-1}A(\xi_1, \dots, \xi_n)) \\ &\leq \mathbf{E}f_{x,\delta}(G^{-1}A(\tilde{\xi}_1, \dots, \tilde{\xi}_n)) + \frac{8}{3} (C_A C_4)^{3/2} C_f \delta^{-3} n^{-1/2}. \end{aligned}$$

It remains to note that

$$\mathbf{E}f_{x,\delta} \left(G^{-1}A(\tilde{\xi}_1, \dots, \tilde{\xi}_n) \right) \leq \mathbf{P} \left(G^{-1}A(\tilde{\xi}_1, \dots, \tilde{\xi}_n) > x - \delta \right)$$

The last statement of the corollary follows from the obvious fact that the density of $G^{-1}A(\tilde{\xi}_1, \dots, \tilde{\xi}_n)$ is bounded by 1 for every G with $G^2 \geq \mathbf{E}|A(\tilde{\xi}_1, \dots, \tilde{\xi}_n)|^2$. \square

6.3. A family of quadratic forms

Here we briefly discuss the situation arising in adaptive testing problem when the maximum of a family of quadratic forms of ξ_i 's is considered. We again aim to show that the joint distribution of this family (and thus the distribution of the maximum) can be well approximated by the similar distribution for quadratic forms of Gaussian random variables.

Let A_1, \dots, A_M be a collection of symmetric $n \times n$ -matrices with vanishing diagonal elements. We analyze the joint distribution of the normalized quadratic forms $G_m^{-1}A_m(\xi_1, \dots, \xi_n)$ with independent random variables ξ_i satisfying $\mathbf{E}\xi_i = 0$, $\mathbf{E}\xi_i^2 = \sigma_i^2$ and $\mathbf{E}\xi_i^4 < \infty$, and some constants G_m , $m = 1, \dots, M$. More precisely, we intend to show that the distribution of this family is close to the distribution of the family $\{G_m^{-1}A_m(\tilde{\xi}_1, \dots, \tilde{\xi}_n), m = 1, \dots, M\}$ with Gaussian variables $\tilde{\xi}_i \sim \mathcal{N}(0, \sigma_i^2)$.

Proposition 6.4. *Let the variables ξ_i fulfill $\mathbf{E}\xi_i^4 \leq C_E\sigma_i^4$ and let every matrix A_m satisfy the conditions of Proposition 6.3 with the same constant C_A , $m = 1, \dots, M$. Then, for every three times continuously differentiable function f in the space \mathbb{R}^M , it holds*

$$\left| \mathbf{E}f \left(G^{-1}A(\xi_1, \dots, \xi_n) \right) - \mathbf{E}f \left(G^{-1}A(\tilde{\xi}_1, \dots, \tilde{\xi}_n) \right) \right| \leq \frac{8}{3}f_3M^3(C_4C_A)^{3/2}n^{-1/2}$$

where $G^{-1}A$ denotes the vector with elements $G_m^{-1}A_m$ and f_3 means the maximum of the absolute value of the third derivative of f , that is,

$$f_3 = \sup_{x \in \mathbb{R}^M} \max_{i,j,k=1,\dots,M} \left| \frac{\partial^3 f(x)}{\partial x_i \partial x_j \partial x_k} \right|.$$

Proof. The proof follows the same line as in the case of one quadratic forms when understanding $G^{-1}A$, u_{n-1} , $f'(u_{n-1})$ and Δ_n as vectors in \mathbb{R}^M and $f''(u_{n-1})$ as the $M \times M$ -matrix of the second derivatives of f at u_{n-1} . The only difference is that we apply the bound $\mathbf{E}|\Delta_n|^3 \leq M^3 8(C_4C_A)^{3/2}n^{-3/2}$ for the norm of Δ_n which is M^3 times larger than in the case of $M = 1$, cf. (6.7). The details are left to the reader. \square

A straightforward corollary of this results concerns the maximum of $G_m^{-1}A_m$'s.

Corollary 6.2. *Let the conditions of Proposition 6.4 be fulfilled. Then*

$$\begin{aligned} \mathbf{P} \left(\max_{m \leq M} G_m^{-1} A_m(\xi_1, \dots, \xi_n) \leq x \right) - \mathbf{P} \left(\max_{m \leq M} G_m^{-1} A_m(\tilde{\xi}_1, \dots, \tilde{\xi}_n) \leq x - \delta \right) \\ \leq \text{Const.} M^3 C_A^{3/2} n^{-1/2} \delta^{-3} \end{aligned}$$

with a constant *Const.* depending on C_4 only. If, in addition, $G_m^2 \geq \mathbf{E}|A_m(\xi_1, \dots, \xi_n)|^2$ for all $m \leq M$, then

$$\begin{aligned} \mathbf{P} \left(\max_{m \leq M} G_m^{-1} A_m(\xi_1, \dots, \xi_n) \leq x \right) - \mathbf{P} \left(\max_{m \leq M} G_m^{-1} A_m(\tilde{\xi}_1, \dots, \tilde{\xi}_n) \leq x \right) \\ \leq \text{Const.} M^3 C_A^{3/2} n^{-1/2} \delta^{-3} + M\delta. \end{aligned}$$

Proof. The first statement can be checked exactly as for the case of $M = 1$, see the proof of Corollary 6.1. As regard to the second statement, it suffices to mention that the density of each $G_m^{-1} A_m(\tilde{\xi}_1, \dots, \tilde{\xi}_n)$ is bounded by 1 and hence the density of the maximum of $G_m^{-1} A_m(\tilde{\xi}_1, \dots, \tilde{\xi}_n)$'s is bounded by M . \square

Remark 6.3. If M is not too large in the sense that $M^3 n^{-1/2}$ is small, then, selecting a proper δ , we can derive from this statement that the distribution of the maximum of $G_m^{-1} A_m(\xi_1, \dots, \xi_n)$'s is approximated by the similar distributions for $G_m^{-1} A_m(\tilde{\xi}_1, \dots, \tilde{\xi}_n)$'s.

References

- [1] Billingsley, P. (1968). *Convergence of probability measures*. J.Wiley. New York.
- [2] Burnashev, M.V. (1979) On the minimax detection of an inaccurately known signal in a white Gaussian noise background. *Theory Probab. Appl.* **24** 107–119.
- [3] Ermakov, M.S. (1990). Minimax detection of a signal in a white Gaussian noise. *Theory Probab. Appl.*, **35**, 667–679.
- [4] Härdle, W. and Mammen, E. (1993). Comparing nonparametric versus parametric regression fits. *Annals of Statistic* **4** 1926–1947.
- [5] Härdle, W., Spokoiny, V. and Sperlich, S. (1995). Semiparametric single index versus fixed link function modeling. Preprint **190**, Weierstrass Institute, Berlin.
- [6] Hart, J. (1997). *Nonparametric Smoothing and Lack-of-Fit Tests* New York, Berlin, Heidelberg: Springer.
- [7] Ibragimov, I.A. and Khasminskii, R.Z. (1977). One problem of statistical estimation in Gaussian white noise. *Soviet Math. Dokl.* **236** no.4, 1351–1354.
- [8] Ingster, Yu.I. (1982). Minimax nonparametric detection of signals in white Gaussian noise *Problems Inform. Transmission* **18** 130–140.
- [9] Ingster, Yu.I. (1993). Asymptotically minimax hypothesis testing for nonparametric alternatives. I–III. *Math. Methods of Statist.* **2** 85–114, **3** 171–189, **4** 249–268.
- [10] Ingster, Yu.I. and Suslina, I. (1998). Minimax nonparametric hypothesis testing for ellipsoids and Besov bodies. *Problems Inform. Transmission* **34** no. 1, 56–68.
- [11] Klein, R.L. and Spady, R.H. (1993). An efficient semiparametric estimator for binary response models. *Econometrica*, **61**, 387–421.

- [12] Lehmann, E.L. (1959) *Testing Statistical Hypothesis* Wiley, New York.
- [13] Lepski, O.V. and Spokoiny, V.G. (1998). Minimax Nonparametric Hypothesis Testing: The Case of an Inhomogeneous Alternative. *Bernoulli*, to appear.
- [14] Lepski, O. and Tsybakov, A. (1996). Asymptotically exact nonparametric hypothesis testing in sup-norm and at a fixed point. SFB Discussion Paper **91**, Humboldt University, Berlin.
- [15] Mann, H.B. and Wald, A.(1942) On the choice of the number of intervals in the application of the chi-square test. *Ann. Math. Stat.* **13** 306–317.
- [16] Neyman, J. (1937). “Smooth test” for goodness of fit. *Scand. Aktuarietidskr.* **20** 149–199.
- [17] Petrov, V.V. (1975). *Sums of Independent Random Variables*. Springer, New York.
- [18] Spokoiny, V. (1996). Adaptive hypothesis testing using wavelets. *Annals of Statistics* **24** no. 6, 2477–2498.
- [19] Spokoiny, V. (1998). Adaptive and spatially adaptive testing a nonparametric hypothesis. *Math. Methods of Stat.* **7**, no. 3, 245–273.
- [20] Triebel, H. (1992). *Theory of function spaces*. Birkhäuser, Basel.

WEIERSTRASS INSTITUTE FOR APPLIED ANALYSIS AND STOCHASTICS, MOHRENSTR. 39, 10117
BERLIN, GERMANY

E-mail address: `spokoiny@wias-berlin.de`