



**HAL**  
open science

# Goodness-of-Fit Tests for Variance Function in Regression Models

Sandie Ferrigno, Marie-José Martinez

► **To cite this version:**

Sandie Ferrigno, Marie-José Martinez. Goodness-of-Fit Tests for Variance Function in Regression Models. CMStatistics 2022, Dec 2022, London, United Kingdom. hal-03935703

**HAL Id: hal-03935703**

**<https://hal.science/hal-03935703>**

Submitted on 12 Jan 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

### OBJECTIVES

Let consider  $(X, Y) \in \mathbb{R}^2$  and a regression model to predict the value of  $Y$  from that of  $X$ ,

$$Y = m(X) + \sigma(X)\varepsilon,$$

where  $m(\cdot)$  is the regression function,  $\sigma^2(\cdot)$  the variance function et  $\varepsilon$  the random error term.

In addition to assumptions about the functional form of the regression function and the variance function, this model requires that the random error term is additive, independent of  $X$  and is often assumed to follow the normal distribution  $N(0, 1)$ . Many tests have been proposed to assess these assumptions. Most of them are “directional” in that

they detect departures from mainly this given assumption of the model. When a number of directional tests are applied to the same model, each requiring the correctness of other assumptions, the assessment of the overall validity of the model may become a difficult matter to untangle.

We focus on the task of choosing the structural part  $\sigma^2(\cdot)$ . Dette *et al.* [1] and Pardo-Fernández *et al.* [2] proposed **directional tests** to valid the form of the variance function while Ducharme *et al.* [3] proposed a **global test**.

We propose to compare these previous tests through simulations.

### MODEL

The simulated model was taken to be

$$Y = 1 + \sin(X) + \sigma(X)\varepsilon \quad (1)$$

where  $\sigma(X) = 0.5 * \exp(cX)$  and the term  $c$ , given by 0, 0.5, 1, represents the deviation from the null hypothesis. The null hypothesis consists in the homoscedastic model  $\sigma(X) = 0.5$  ( $c = 0$ ). We also suppose that  $X \sim U([0; 1])$ ,  $\varepsilon \sim N(0, 1)$ , and  $\varepsilon$  and  $X$  are independent. We use **Wild bootstrap** [4] procedures to calculate the critical test value for finite samples. The null hypothesis is rejected if the test statistic is bigger than the  $(1 - \alpha)$ -quantile ( $\alpha = 0.05$ ) of the wild bootstrap distribution of the test statistic.

### RESULTS AND CONCLUSIONS

For each combination of factors, the experiment is replicated 50 times. The wild bootstrap resampling has been performed 50 times for each sample. We use the Epanechnikov kernel and optimal bandwidths to calculate nonparametric estimators of the different functions used in the three tests. We compare the results for two sample sizes :  $n = 100$  and  $n = 200$ .

$n = 100$	Test1	Test2	Test3
$c = 0$	0.02	0.09	0.06
$c = 0.5$	0.74	0.92	0.24
$c = 1$	0.96	1	0.66

$n = 200$	Test1	Test2	Test3
$c = 0$	0.04	0.06	0.03
$c = 0.5$	0.92	0.96	0.5
$c = 1$	1	1	0.72

i) Under the null hypothesis, the proportions of rejection are similar to the theoretical level, especially for the two directional tests and  $n = 200$ .  
ii) Under the alternatives, the different results are globally better as  $n$  increases but the global test is not as efficient as the directional tests.

More complete results will be given in a forthcoming article and a tool will be offered to users of this type of model to test their validity in different ways.

### DETTE *et al.* (TEST1)

**Test hypothesis:**

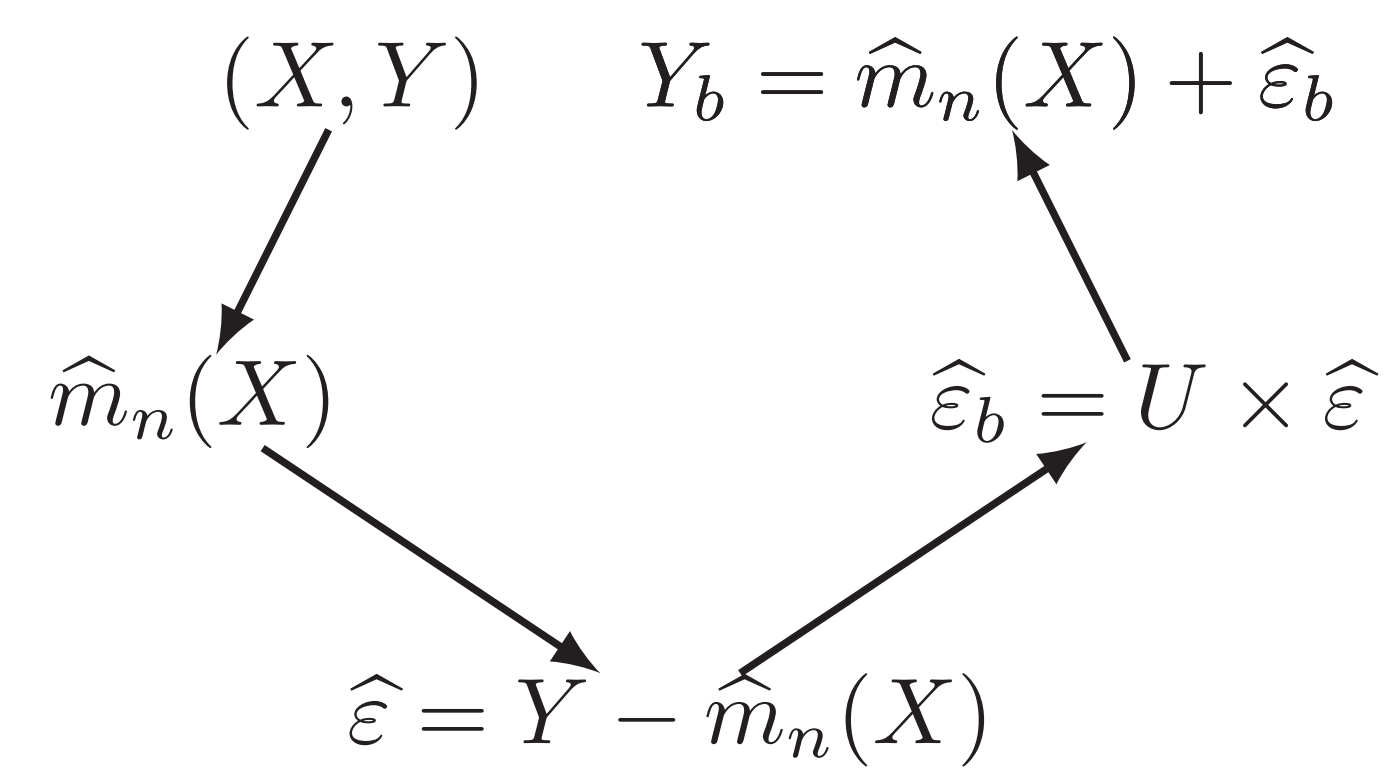
$$H_0 : \sigma^2(x) = \sigma_0^2(x) \quad \text{vs} \quad H_1 : \sigma^2(x) \neq \sigma_0^2(x)$$

where  $\sigma_0^2(x)$  is the variance function associated to the model (1). The test statistic (under  $H_0$ ) is:

$$T_n = n \int \left( \widehat{F}_\varepsilon(x) - \widehat{F}_{\varepsilon_0}(x) \right)^2 d\widehat{F}_\varepsilon(x)$$

where  $\varepsilon = \{Y - m(X)\}/\sigma(X)$  are the standardized residuals and  $\widehat{F}_\varepsilon(x)$  is the Nadaraya-Watson nonparametric estimate of the distribution function  $F_\varepsilon$ .

**Wild bootstrap procedure and test statistic:**



$$T_n^b = n \int \left( \widehat{F}_\varepsilon(y) - \widehat{F}_\varepsilon^b(y) \right)^2 d\widehat{F}_\varepsilon^b(y)$$

where  $U$  is Normal, Mammen or Rademacher distribution and  $\widehat{F}_\varepsilon^b(y)$  is calculated from  $(X, Y_b)$ .

### PARDO-FERNÁNDEZ *et al.* (TEST2)

**Test hypothesis:**

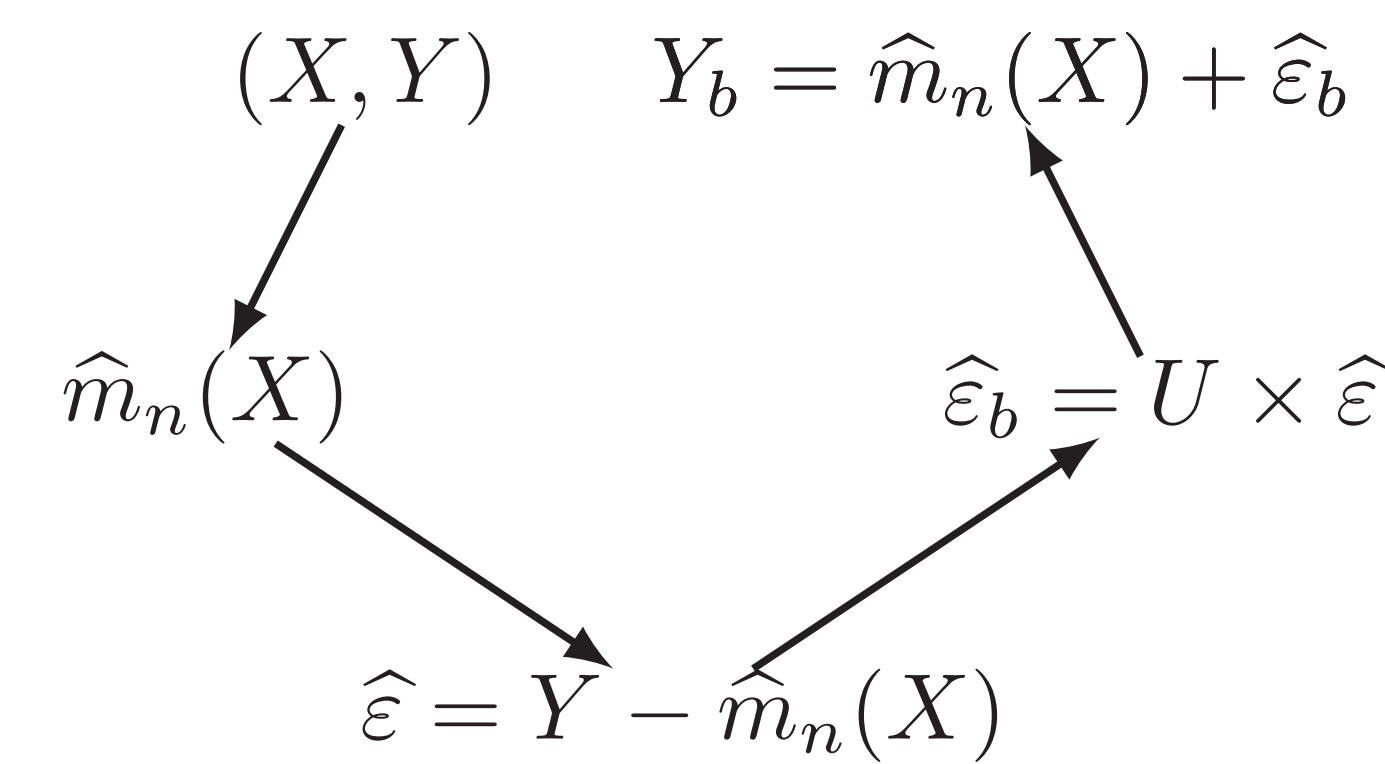
$$H_0 : \sigma^2(x) = \sigma_0^2(x) \quad \text{vs} \quad H_1 : \sigma^2(x) \neq \sigma_0^2(x)$$

where  $\sigma_0^2(x)$  is the variance function associated to the model (1). The test statistic (under  $H_0$ ) is:

$$T_n = n \int \left( \widehat{\phi}(y) - \phi_0(y) \right)^2 w(y) dy$$

where  $\widehat{\phi}(y)$  is the Nadaraya-Watson nonparametric estimator of the characteristic function and  $w(\cdot)$  the pdf of a normal distribution.

**Wild bootstrap procedure and test statistic:**



$$T_n^b = n \int \left( \widehat{\phi}(y) - \widehat{\phi}^b(y) \right)^2 w(y) dy$$

where  $U$  is Normal, Mammen or Rademacher distribution and  $\widehat{\phi}^b(y)$  is calculated from  $(X, Y_b)$ .

### DUCHARME *et al.* (TEST3)

**Test hypothesis:**

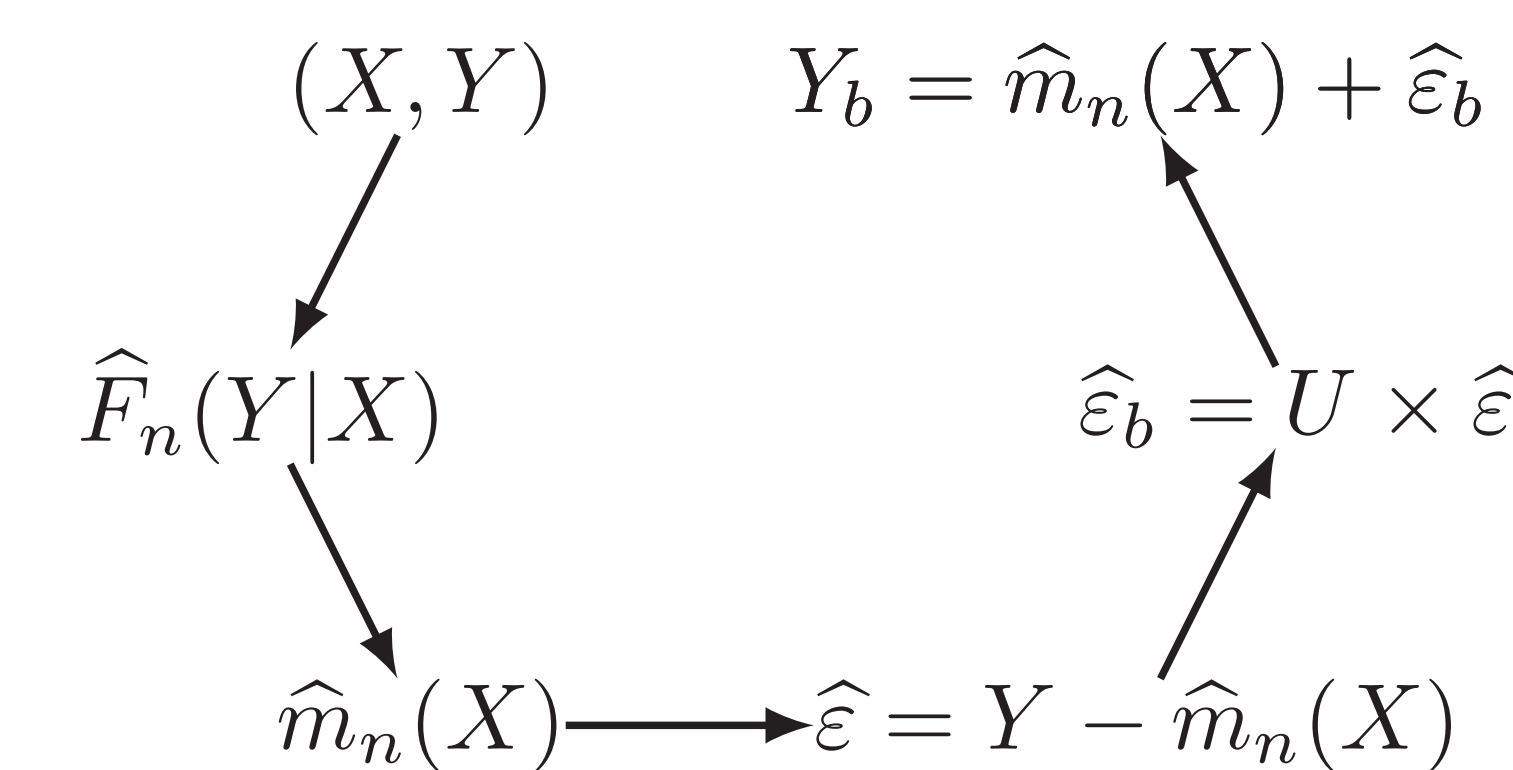
$$H_0 : F(y|x) = F_0(y|x) \quad \text{vs} \quad H_1 : F(y|x) \neq F_0(y|x)$$

where  $F_0(y|x)$  is the conditionnal distribution associated to the model (1). The test statistic (under  $H_0$ ) is:

$$T_n = n\sqrt{h} \int \int \left( \widehat{F}_n(y|x) - F_0(y|x) \right)^2 F_0(dy|x) dx$$

where  $\widehat{F}_n(y|x)$  is the local linear estimator of the conditional distribution function and  $h$  a bandwidth.

**Wild bootstrap procedure and test statistic:**



$$T_n^b = n\sqrt{h} \int \int \left( \widehat{F}_n(y|x) - \widehat{F}_n^b(y|x) \right)^2 \widehat{F}_n^b(dy|x) dx$$

where  $U$  is Normal, Mammen or Rademacher distribution and  $\widehat{F}_n^b(y|x)$  is calculated from  $(X, Y_b)$ .

### REFERENCES

- [1] Dette, Neumeyer, and Van Keilegom. A new test for the parametric form of the variance function in non-parametric regression. *J.R.Statist.Soc.B*, 69:903–917, 2007.
- [2] Pardo-Fernández and Jiménez-Gamero. A model specification test for the variance function in non-parametric regression. *ASTA Advances in Statistical analysis*, 103:387–410, 2019.
- [3] Ducharme and Ferrigno. An omnibus test of goodness-of-fit for conditional distributions with applications to regression models. *Journal of Statistical Planning and Inference*, 142:2748–2761, 2012.
- [4] Davidson and Flachaire. The wild bootstrap, tamed at last. *Journal of Econometrics*, 146:162–169, 2008.