

## Comments on: An updated review of Goodness-of-Fit tests for regression models

Ingrid Van Keilegom

Published online: 25 July 2013

© Sociedad de Estadística e Investigación Operativa 2013

First of all I would like to congratulate the authors for writing this very interesting, timely and thorough review paper. It is an almost impossible task to review the very extensive literature on smoothing-based goodness-of-fit tests in regression in one single paper. Nevertheless, the paper contains a very rich collection of the contributions over the last 10–20 years, not only about goodness-of-fit tests for the regression function, but also regarding the variance function, the error distribution, the case of incomplete or dependent data, and much more. A comprehensive book on goodness-of-fit tests in regression using smoothing-based approaches is somewhat missing in the literature, and I hope that the authors will use this paper as a starting point for a more extensive book on this topic.

In this discussion I would like to comment on three issues, which the authors did not consider in their review: (1) Goodness-of-fit tests for the coefficient of variation; (2) tests for the independence between the error term and the covariates; and (3) some open questions and comments.

### 1 Goodness-of-fit tests for the coefficient of variation

The coefficient of variation for  $X = x$  is given by the ratio  $m(x)/\sigma(x)$  with  $m(x) = E(Y|X = x)$  and  $\sigma^2(x) = \text{Var}(Y|X = x)$ , and is an important quantity in regression. For instance, numerous authors have considered the problem of estimating the regression function assuming that the coefficient of variation is constant as a function of  $x$

---

This comment refers to the invited paper available at doi:[10.1007/s11749-013-0327-5](https://doi.org/10.1007/s11749-013-0327-5).

I. Van Keilegom (✉)

Université catholique de Louvain, Louvain-la-Neuve, Belgium

e-mail: [ingrid.vankeilegom@uclouvain.be](mailto:ingrid.vankeilegom@uclouvain.be)

(see e.g. Eagleson and Müller 1997, among others). Other authors, like e.g. Engle et al. (1987), were interested in the relation between  $m(\cdot)$  and  $\sigma(\cdot)$  in econometric and financial models, since the return (or the regression function,  $m$ ) and the risk (or the scale function,  $\sigma$ ) are expected to be related. Finally, many time series models can be written as a regression model of the form  $Y = m(X) + \sigma(X)\varepsilon$  with  $m(X) = c\sigma(X)$  for some constant  $c$ . This is the case for e.g. the ARCH model, and any other time series model with a multiplicative structure of the form  $Z_t = \sigma_t \varepsilon_t$ , where  $\varepsilon_t$  has mean 0 and variance 1 and is independent of  $Z_{t-1}$ , and where  $\sigma_t = g(Z_{t-1})$  for some function  $g$ . Hence in order to test whether the time series has a multiplicative structure one can test whether the corresponding regression model has a constant coefficient of variation.

In the recent literature on non-parametric regression a number of tests have been proposed for the proportionality of  $m(\cdot)$  and  $\sigma(\cdot)$ . Dette et al. (2009) proposed a test based on the comparison between an estimator of the error distribution under the null hypothesis of proportionality and an estimator that does not make this assumption. Their test is similar in spirit to the test for the regression function proposed in Sect. 2.4. Dette and Wieczorek (2009) proposed to test for proportionality of  $m(\cdot)$  and  $\sigma(\cdot)$  by looking at the difference between the square of the regression function and the product of the proportionality constant and the variance function. Their test statistic is a weighted  $L_2$ -distance between non-parametric kernel estimators of these two quantities. Recently, Dette et al. (2012) proposed a test that is based on a comparison between the ‘observations’  $Y_i/\hat{\sigma}(X_i)$  ( $i = 1, \dots, n$ ) and their conditional mean (which is constant under the hypothesis of proportionality), where  $\hat{\sigma}(\cdot)$  is a local linear estimator of the scale function  $\sigma(\cdot)$ .

## 2 Tests for the independence between the error term and the covariates

Consider the non-parametric location-scale model

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

where  $\varepsilon$  and  $X$  are independent. Numerous goodness-of-fit tests for  $m(\cdot)$  and  $\sigma(\cdot)$  have been constructed assuming that  $\varepsilon$  and  $X$  are independent (see Sect. 2.4 and subsequent sections for more details). It is therefore important to have at hand powerful testing procedures for this independence assumption. Einmahl and Van Keilegom (2008a) constructed Kolmogorov–Smirnov, Cramér–von Mises and Anderson–Darling type test statistics based on a non-parametric kernel estimator of the distance between  $F_{X,\varepsilon}(x, e)$  and  $F_X(x)F_\varepsilon(e)$  for all  $(x, e)$  in the support of  $(X, \varepsilon)$ . They use a smoothed bootstrap to calibrate their test.

Einmahl and Van Keilegom (2008b) use another approach for the special case of a homoscedastic model, i.e. when  $\sigma(\cdot) \equiv \sigma$ . They argue that when the function  $m$  is smooth, then  $Y_{[i+1]} - 2Y_{[i]} + Y_{[i-1]}$  is approximately equal to  $\varepsilon_{[i+1]} - 2\varepsilon_{[i]} + \varepsilon_{[i-1]}$ . Here,  $Y_{[1]}, \dots, Y_{[n]}$  are the concomitants, i.e. the  $Y$ -values corresponding to the ordered  $X$ -values  $X_{(1)}, \dots, X_{(n)}$ , and similarly for  $\varepsilon_{[1]}, \dots, \varepsilon_{[n]}$ . They then check the independence between  $\varepsilon$  and  $X$  by verifying whether  $Y_{[i+1]} - 2Y_{[i]} + Y_{[i-1]}$  and  $X_{(i)}$  are independent. This approach has the important advantage that  $m(\cdot)$  does not need to be estimated, avoiding hence the delicate choice of a smoothing parameter.

Neumeyer (2009) proposed to test the independence between  $\varepsilon$  and  $X$  by looking at the difference between the marginal distribution of  $X$  and the conditional distribution of  $X$  given  $\varepsilon$ . She proposes a test statistic that is a weighted kernal based  $L_2$ -distance between estimators of these two distributions. The so-obtained test statistic is a  $U$ -statistic, for which she shows the asymptotic normality. She also proves how to obtain a valid bootstrap procedure to approximate the distribution of this test statistic.

Finally, Hlávka et al. (2011) proposed a test based on the Fourier formulation of independence, and they utilize the joint and the marginal empirical characteristic functions of  $X$  and of a non-parametric kernel estimator of  $\varepsilon$  to test independence. The spirit of their test is rather similar to the idea of the tests based on the characteristic function that were proposed by Hušková and Meintanis (2007, 2009, 2010) in other contexts (see Sect. 2.4 for more details).

### 3 Some open questions and comments

Although the literature on goodness-of-fit tests has known an enormous expansion over the last 10 years and many testing problems have been successfully studied, there are still a number of open problems that need closer attention. Below I give three examples of areas in which goodness-of-fit tests are not much or not at all developed. There are certainly many more examples of areas in which goodness-of-fit tests need to be developed. The examples below are selected merely based on personal interests and expertise.

First of all, it would be interesting to consider test statistics for the hypothesis which states a certain relationship between  $m(\cdot)$  and  $\sigma(\cdot)$ . In Sect. 1 of this discussion I mentioned tests for proportionality of  $m(\cdot)$  and  $\sigma(\cdot)$  that have been proposed in the literature on non-parametric regression, but other non-constant relationships are also of interest in practice. In a paper in preparation, Escanciano et al. (2013) are developing a test for this hypothesis, but more research in this direction is needed.

Another goodness-of-fit problem that has not received much attention in the literature is the goodness-of-fit of general semi- and non-parametric transformation models of the form

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

where  $\varepsilon$  and  $X$  are independent (for identifiability reasons), and the transformation  $\Lambda$  is assumed to belong to a family of parametric transformations (like e.g. the famous Box–Cox family) or  $\Lambda$  can be estimated non-parametrically. The identifiability and estimation of this model has received quite some attention in the literature (see e.g. Linton et al. 2008, and the references therein), but to the best of my knowledge no-one has considered goodness-of-fit tests for the regression function under this model, or more generally tests for the appropriateness of the transformation model itself.

A final example of an area in which goodness-of-fit tests seem to be undeveloped so far is the area of cure models in survival analysis. When modeling time-to-event data (like e.g. the time to death of a patient due to a certain disease), we typically assume that all subjects are at risk and will experience the event of interest if followed long enough. However, a typical feature of many medical applications is the possibility of ‘cure’, in the sense that some of the subjects will actually not experience the

event. Cure models are survival models allowing a cured proportion of individuals. Moreover, measuring times to a certain event in practice naturally induces the presence of right censoring. Combining both censoring and possibility of cure involves identifiability problems. Indeed, even though the follow-up period is long, it is hard to distinguish a censored individual in the uncured group from a cured individual. In the literature there are two major classes of cure models: the mixture cure model (see e.g. the book by Maller and Zhou 1996, and the references therein) and the promotion time cure model (see e.g. Tsodikov 1998). To the best of my knowledge, formal goodness-of-fit procedures for testing the appropriateness of these two models have not been developed so far.

I would be interested to hear the authors thoughts on the above raised issues, and on potential approaches based on the literature they surveyed.

**Acknowledgements** I. Van Keilegom acknowledges financial support from the European Research Council under the European Community's Seventh Framework Programme (FP7/2007-2013)/ERC Grant agreement No. 203650, from IAP research network P7/06 of the Belgian Government (Belgian Science Policy), and from the contract 'Projet d' Actions de Recherche Concertées' (ARC) 11/16-039 of the 'Communauté française de Belgique', granted by the 'Académie universitaire Louvain'.

## References

- Dette H, Wiecek G (2009) Testing for a constant coefficient of variation in nonparametric regression. *J Stat Theory Pract* 3:587–612
- Dette H, Pardo-Fernández JC, Van Keilegom I (2009) Goodness-of-fit tests for multiplicative models with dependent data. *Scand J Stat* 36:782–799
- Dette H, Marchlewski M, Wagener J (2012) Testing for a constant coefficient of variation in nonparametric regression. *Ann Inst Stat Math* 64:1045–1070
- Eagleson GK, Müller HG (1997) Transformations for smooth regression models with multiplicative errors. *J R Stat Soc, Ser B* 59:173–189
- Einmahl J, Van Keilegom I (2008a) Specification tests in nonparametric regression. *J Econom* 143:88–102
- Einmahl J, Van Keilegom I (2008b) Tests for independence in nonparametric regression. *Stat Sin* 18:601–616
- Engle RF, Lilien DM, Robins RP (1987) Estimating time varying risk premia in the term structure: the Arch-M model. *Econometrica* 55:391–407
- Escanciano JC, Pardo-Fernández JP, Van Keilegom I (2013) A nonparametric test for risk-return relationships (in preparation)
- Hlávka Z, Hušková M, Meintanis SG (2011) Tests for independence in non-parametric heteroscedastic regression models. *J Multivar Anal* 102:816–827
- Hušková M, Meintanis S (2007) Omnibus tests for the error distribution in linear regression models. *Statistics* 41:363–376
- Hušková M, Meintanis S (2009) Goodness-of-fit tests for parametric regression models based on empirical characteristic functions. *Kybernetika* 45:960–971
- Hušková M, Meintanis S (2010) Test for the error distribution in nonparametric possibly heteroscedastic regression models. *Test* 19:92–112
- Linton O, Sperlich S, Van Keilegom I (2008) Estimation of a semiparametric transformation model. *Ann Stat* 36:686–718
- Maller R, Zhou X (1996) *Survival analysis with long-term survivors*. Wiley, New York
- Neumeyer N (2009) Testing independence in nonparametric regression. *J Multivar Anal* 100:1551–1566
- Tsodikov A (1998) A proportional hazards model taking account of long-term survivors. *Biometrics* 54:1508–1516