

## **ON HETEROSCEDASTICITY IN ROBUST REGRESSION**

**Jan Kalina**

---

### **Abstract**

This work studies the phenomenon of heteroscedasticity and its consequences for various methods of linear regression, including the least squares, least weighted squares and regression quantiles. We focus on hypothesis tests for these regression methods. The new approach consists in deriving asymptotic heteroscedasticity tests for robust regression, which are asymptotically equivalent to standard tests computed for the least squares regression. One approach to modeling heteroscedasticity assumes a prior knowledge or specific model for the variability of random regression errors. Another (and more general) approach does not assume a specific form of heteroscedasticity. The paper also describes heteroscedastic regression, which is a tool to incorporate heteroscedasticity to the model. This allows us to define the heteroscedastic least weighted squares regression.

**Key words:** robust statistics, linear regression, diagnostics

**JEL Code:** C14, C12, C21

---

### **Introduction**

Homoscedasticity is one of essential assumptions of linear regression not only for the least squares estimator of regression parameters, but also for its robust counterparts. The paper starts by describing heteroscedasticity as the violation of homoscedasticity and presents its negative consequences. Tests of heteroscedasticity are presented in Section 2 for the least squares estimator, namely the tests of Goldfeld-Quandt, Breusch-Pagan and White test. The new result is the asymptotic version of these tests derived for some robust regression estimator: the least weighted squares (Section 3) and regression quantiles (Section 4). The solution of estimating parameters in the heteroscedastic model is called heteroscedastic regression, which is described again for various regression estimators in Section 5.

### **1 Linear regression**

In the whole paper we consider the linear regression model

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + e_i, i = 1, 2, \dots, n. \quad (1)$$

The variance of the disturbances  $\sigma^2$  is known to be a nuisance parameter. This paper however has also the aim to show that  $\sigma^2$  is a key parameter also in estimating  $\beta$ . It is crucial to estimate  $\sigma^2$  reliably in order to obtain a reliable tests of hypotheses about  $\beta$  and also its reliable confidence intervals. The homoscedasticity assumption

$$\text{var } e_i = \sigma^2, \quad i=1, \dots, n, \quad (2)$$

is called homoscedasticity, while its violation is denoted as heteroscedasticity.

There can be severe negative consequences of heteroscedasticity, especially if the equality of variances of the disturbances is violated heavily. Regression parameters  $\beta$  cannot be estimated efficiently. Denoting the least squares estimator of  $\beta$  by  $b$ , the classical estimator of  $\text{var } b$  is biased. This disqualifies using classical hypothesis tests and confidence intervals for  $\beta$  as well as the value of the coefficient of determination  $R^2$ . Diagnostic tools checking the assumption of equality of variances of the disturbances can be based on residuals  $u = (u_1, \dots, u_n)^T$ , where

$$u_i = Y_i - b_0 - b_1 x_{i1} - \dots - b_p x_{ip} \quad (3)$$

and  $^T$  denotes a vector transposition. In Section 2 we describe objective diagnostic tests.

This paper is devoted entirely to the linear regression model. While in Kalina (2010) we have studied robust statistical methods for the model of multivariate location and scatter, in this paper we point out that also robust multivariate methods are sensitive to the assumption of homoscedasticity. Therefore we must interpret correctly that the variance of disturbances is a nuisance parameter. Particularly in the linear regression it is known that  $\sigma^2$  is a nuisance parameter in estimating the regression parameters  $\beta$ . This does not mean that  $\sigma^2$  is not important or that its estimation stands aside during the inference of  $\beta$ . We bring arguments that  $\sigma^2$  plays a very important role in the statistical inference and influences the estimation procedures, which aim only at the regression parameters  $\beta$ . While the regression is based on the (very non-robust) sum of squares of residuals, the estimation of  $\sigma^2$  is based exactly on the same sum of squares. This connects the problem of non-robustness of estimating  $\beta$  and  $\sigma^2$ .

## **2 Heteroscedasticity for least squares**

We describe the classical Goldfeld-Quandt test, Breusch-Pagan test and White test for the least squares regression. Each of these tests is designed for a different alternative hypothesis. Therefore we do not attempt to summarize the basic context common to all three methods. More details on standard heteroscedasticity tests can be found in econometric references

(Greene, 2002) or (Judge et al., 1985). We point out that the tests were originally proposed in econometric journals, while they are diagnostic tools for a general statistical (not only econometric) context.

Goldfeld-Quandt test (Goldfeld and Quandt, 1965) is easy to be computed and interpreted. It tests the null hypothesis

$$H_0: \text{var } e_i = \sigma^2, \quad i=1, \dots, n, \quad (4)$$

against the alternative hypothesis

$$H_1: \text{var } e = \sigma^2 \text{diag}\{k_1, \dots, k_n\}, \quad i=1, \dots, n, \quad (5)$$

which models heteroscedasticity in a particular way. The constants  $k_1, \dots, k_n$  must be selected by the statistician already before the computation. In fact the test does not depend on these values, but its power depends on them. The alternative hypothesis expresses that the variance of the disturbances depends on some variable (or a combination of variables) in a monotone way. Typically one of the regressors in the linear regression model or fitted values of the response are selected to explain the variability of the disturbances in this way. The test is based on dividing the data to three groups according the values of the constants  $k_1, \dots, k_n$ . Let  $SSE_1$  denote the residual sum of squares in the first group of the data and let  $SSE_3$  denote the residual sum of squares computed in the third group. Let  $r_1$  denote the number of observations in the first group,  $r_3$  in the third group and  $p$  is the number of regression parameters in the linear regression model. The test statistic

$$F = \frac{SSE_1}{SSE_3} \frac{r_1 - p}{r_3 - p} \quad (6)$$

follows Fisher's  $F$ -distribution with  $r_3 - p$  and  $r_1 - p$  degrees of freedom.

Breusch-Pagan test (Breusch and Pagan, 1979) requires to specify the alternative hypothesis of heteroscedasticity in the form

$$\text{var } e_i = \alpha_0 + \alpha_1 Z_{1i} + \dots + \alpha_K Z_{Ki}, \quad i=1, \dots, n, \quad (7)$$

for some variables

$$Z_1 = (Z_{11}, \dots, Z_{1n})^T, \dots, Z_K = (Z_{K1}, \dots, Z_{Kn})^T. \quad (8)$$

Often one or more regressors in the original linear regression model are selected as these auxiliary variables. The null hypothesis corresponds to

$$H_0: \alpha_1 = \alpha_2 = \dots = \alpha_K = 0, \quad (9)$$

which is tested against a general alternative hypothesis that the null hypothesis is not true. Breusch and Pagan (1979) derived the test statistic in the form of a so-called score test, which is one of general asymptotic tests based on the likelihood function, in our case under the presence of nuisance parameters. This tests assumes normal distribution of the disturbances  $e$ .

Víšek (2001) proposed a general test which is known as White test. The test exploits White's proposal of an estimator of the variance matrix  $var e$ , which is consistent also under heteroscedasticity. The test is based on comparing two estimators of the variance matrix, where the classical estimator is consistent only under homoscedasticity, while the White's estimator is consistent also under the alternative hypothesis. Therefore large values of the test statistic speak in favour of the alternative hypothesis.

We would like to point out that White test is a special case of Breusch-Pagan test. Here the particular choice of auxiliary variables  $Z_1, \dots, Z_K$  is performed to contain squares of all regressors in the original model and also products of pairs of regressors in the form  $X_i X_j$  for  $i \neq j$ .

The least squares estimator is known to be too vulnerable with respect to violation of the assumption of the normal distribution of the disturbances  $e$ . Therefore robust statistical methods are studied (see Jurečková and Sen, 1996), which represent a diagnostic tool for the least squares estimator or they can be used as an independent tool for the statistical modeling. One of efficient estimator is the least weighted squares proposed by Vášek (2004), which will now presented.

### **3 Heteroscedasticity for least weighted squares**

We recall the definition of the least weighted squares (LWS) regression estimator and describe asymptotic heteroscedasticity tests, which can be used as diagnostic tools for the LWS regression. The tests are based on the test statistics of the Goldfeld-Quandt, Breusch-Pagan or White test computed for residuals of the least weighted squares.

The least weighted squares (LWS) regression is a robust regression method with a high breakdown point proposed by Vášek (2004). There must be nonnegative weights  $w_1, w_2, \dots, w_n$  specified before the computation of the estimator. While the classical weighted regression assigns a fixed and known weight to each observation, in the context of least weighted squares only the magnitudes of the weights are known a priori. These are assigned to the data after a permutation, which is determined automatically only during the computation based on

the residuals. It is reasonable to choose such weights so that the sequence  $w_1, w_2, \dots, w_n$  is decreasing (non-increasing), so that the most reliable observations obtain the largest weights, while outliers with large values of the residuals get small (or zero) weights.

Let us denote the  $i^{\text{th}}$  order value among the squared residuals for a particular value of the estimate  $\mathbf{b}$  of the parameter  $\boldsymbol{\beta}$  by  $u_i^2(\mathbf{b})$ . The least weighted squares estimator  $\mathbf{b}_{LWS}$  for the model (1) is defined as

$$\mathbf{b}_{LWS} = \operatorname{argmin} \sum_{i=1}^h w_i u_{(i)}^2(\mathbf{b}). \quad (10)$$

Kalina (2007) proposed an approximative algorithm for the intensive computation of the LWS estimator and described diagnostic tests for the estimator, which are equivalent with those computed for the least squares regression.

The least weighted squares estimator has interesting applications, which follow from its robustness and at the same time efficiency for normal data. Theoretical properties including the breakdown point of the estimator are studied by Víšek (2004). It is especially suitable to use the LWS estimator rather than other robust regression estimators, because diagnostic tools (such as tests of heteroscedasticity and autocorrelation of the errors  $\mathbf{e}$ ) can be computed directly using the weighted residuals and again are not affected by outliers. Another advantage of the estimator is that no detection of outliers is actually needed to compute it, because outlying data are downweighted automatically. Víšek (2010) conjectures that the LWS estimator is a reasonable compromise between the least squares and least trimmed squares, namely the estimator combines the efficiency of the least squares with the robustness of the least trimmed squares.

We give an overview of recent results on heteroscedasticity tests for robust regression. Kalina (2009) proposed the asymptotic Goldfeld-Quandt test and the asymptotic Breusch-Pagan test for the least weighted squares estimator. Plát (2004) described the White test for the least weighted squares regression; the paper studies the test statistic of White test computed with the LWS residuals and derives the asymptotic properties of the statistic under the null hypothesis of homoscedasticity. Víšek (2010) derives a more general result of White's estimator of  $\operatorname{var} e$ , which is based on the LWS estimation and is consistent under heteroscedasticity. This allows to define directly a test statistic of White (1980), which is tailor-made for the context of the LWS regression. Now we use these existing results and the ideas of proofs to derive asymptotic heteroscedasticity tests for regression quantiles.

#### 4 Heteroscedasticity for regression quantiles

Regression quantiles represent a natural generalization of sample quantiles to the linear regression model. Their theory is studied by Koenker (2005) and their asymptotic representation was derived by Jurečková and Sen (1996). The estimator depends on a parameter  $\alpha$  in the interval  $(0,1)$ , which corresponds to dividing the disturbances to  $\alpha \cdot 100\%$  values below the regression quantile and the remaining  $(1-\alpha) \cdot 100\%$  values above the regression quantile. Here we describe asymptotic heteroscedasticity tests for regression quantiles, which are derived based on their asymptotic representation. The proof of the theorems follows from the asymptotic considerations of Kalina (2009).

*Theorem.* Let the test statistic  $F$  of the Goldfeld-Quandt test be computed using residuals of the quantile regression estimator with a parameter  $\alpha$ . Then  $F$  has asymptotically Fisher's  $F$ -distribution with  $r_3 - p$  and  $r_1 - p$  degrees of freedom under the null hypothesis of homoscedasticity and assuming normal distribution of disturbances in the linear regression model.

*Theorem.* Let the regression quantile estimator with parameter  $\alpha$  be computed in the linear regression model. Let the test statistic of Breusch-Pagan test  $\chi^2$  be computed as one half of regression sum of squares in the model

$$\frac{u_i^2}{s^2} = \alpha_0 + \alpha_1 Z_{1i} + \dots + \alpha_K Z_{Ki} + v_i, \quad i=1, \dots, n, \quad (11)$$

where  $u = (u_1, \dots, u_n)^T$  is the vector of residuals of the regression quantile estimator and  $s^2$  is the estimator of  $\sigma^2$ . Then the test statistic  $\chi^2$  is asymptotically  $\chi_K^2$  distributed assuming the null hypothesis of homoscedasticity and normal distribution of disturbances in the linear regression model.

#### 5 Heteroscedastic least weighted squares regression

If the null hypothesis of equality of variances in the model (1) is rejected by one of the tests of Section 2 (for least squares), Section 3 (for least weighted squares) or Section 4 (for regression quantiles), we recommend to transform the model (1) to another model in order to suppress the negative consequences of heteroscedasticity. The estimation of regression

parameters in the transformed model is called heteroscedastic regression. We discuss the procedure on the example of the LWS regression.

Assumptions or a prior knowledge on the form of heteroscedasticity should be incorporated within the process of removal heteroscedasticity. This is the case of the Goldfeld-Quandt test in the formula (5). Using the same notation we work with the model

$$\frac{Y_i}{\sqrt{k_i}} = \frac{\beta_1 X_{1i}}{\sqrt{k_i}} + \dots + \frac{\beta_p X_{pi}}{\sqrt{k_i}} + \frac{e_i}{\sqrt{k_i}}, \quad i=1, \dots, n. \quad (12)$$

One of typical examples is the choice  $\sqrt{k_i} = X_{ji}$  for a certain  $j$  and for  $i=1, \dots, n$ , where the variance of the errors is modeled to be directly proportional to the  $j$ -th regressor. Other examples include  $\sqrt{k_i} = \sqrt{X_{ji}}$  or  $\sqrt{k_i} = \hat{Y}_i = b_1 X_{1i} + \dots + b_p X_{pi}$ , where  $i=1, \dots, n$ . In the model (12) we estimate the regression parameters by the least weighted squares method and heteroscedasticity should be tested again. If the null hypothesis of homoscedasticity is not rejected this time, then the model (12) is considered to be preferable to the model (1). Therefore we consider only the results of the transformed model (12) including not only the point estimates of  $\beta$ , but also confidence intervals and hypothesis tests of  $\beta$  based on the asymptotic distribution of the LWS estimator, the value of the robust coefficient of determination and other statistics.

However sometimes the variability of the disturbances is modeled in a more complicated way, just like in formula (7) in the Breusch-Pagan test. Then we describe a possible procedure for the removal of heteroscedasticity in two stages. In the *first stage* the regression parameters in the model (1) are estimated by the least weighted squares method and squares of the LWS residuals  $u_i^2$  are computed. Then the regression parameters in the auxiliary regression model

$$u_i^2 = \alpha_0 + \alpha_1 Z_{1i} + \dots + \alpha_K Z_{Ki} + v_i, \quad i=1, \dots, n, \quad (13)$$

are estimated by the LWS estimator, where  $v_1, \dots, v_n$  are random disturbances. Thus we obtain estimates  $\hat{\alpha}_0, \hat{\alpha}_1, \dots, \hat{\alpha}_K$  for regression parameters  $\alpha_0, \alpha_1, \dots, \alpha_K$ . In the *second stage* the fitted values of  $u_i^2$ , which are computed as

$$\hat{u}_i^2 = \hat{\alpha}_0 + \hat{\alpha}_1 Z_{1i} + \dots + \hat{\alpha}_K Z_{Ki}, \quad i=1, \dots, n, \quad (14)$$

are used as the constants  $k_1, \dots, k_n$  for the transformed model (12), in which the estimators are computed using the LWS estimation procedure.

White test is often understood as a general method, which does not contain any recommendation about a possible removal of the heteroscedasticity (Greene, 2002). However since it is a special case of the Breusch-Pagan test, it also allows the heteroscedastic regression to be used in the same spirit. Therefore if the White test as a diagnostic tool for the LWS regression gives a significant result, the heteroscedastic regression (13)-(14) should be applied and the squares of all regressors and products of pairs of regressors (Section 2) are a natural choice for the auxiliary variables for the regression model (13), in which the LWS regression can be used to estimate the regression parameters.

## **Conclusion**

This work studies the phenomenon of heteroscedasticity in robust regression. Assuming the standard linear regression model, the consequences of heteroscedasticity for robust regression are described and asymptotic heteroscedasticity tests for the least weighted squares regression and for regression quantiles are derived.

We also describe two possible ways of removing heteroscedasticity from the linear regression model. Both are based on a transformation of the original model and take into account such variables, which could possibly explain the variability of the disturbances. In other words this approach models the heteroscedasticity in a particular way. In practice such modeling is based on prior assumptions or knowledge. There exists no heteroscedasticity test optimal uniformly over all situations, but rather different tests have different properties. Therefore it is not possible to select the optimal heteroscedasticity test for a given data set. Another possibility is to use a robust regression estimator consistent also under the assumption of heteroscedastic disturbances (Víšek, 2010). Nevertheless our approaches may be more appropriate for high-dimensional data.

## **Acknowledgment**

This research is supported by the grant 402/09/0557 (Robustification of selected econometric methods) of the Grant Agency of the Czech Republic.

## **References**

1. Breusch, T.S., Pagan A.R. Simple test for heteroscedasticity and random coefficient variation. *Econometrica* 47 (5), 1287-1294, 1979.
2. Goldfeld, S.M., Quandt R.E. Some tests for homoscedasticity. *Journal of the American Statistical Association* 60 (310), 539-547, 1965.



3. Greene, W.H. *Econometric analysis*. Fifth edition. New York: Macmillan, 2002.
4. Judge, G.G., Griffiths, W.E., Hill, R.C., Lutkepohl, H., Lee, T.C. *The theory and practice of econometrics*. New York: Wiley, 1985.
5. Jurečková, J., Sen, P.K. *Robust statistical procedures: Asymptotics and interrelations*. New York: Wiley, 1996.
6. Kalina, J. Implicitly weighted multivariate methods for high-dimensional data. In *MSED 2010 – Mezinárodní statisticko-ekonomické dny na VŠE: Sborník příspěvků. Praha, 9.-10. 9. 2010*. University of Economics, Prague, CD-ROM, 1-7, 2010.
7. Kalina, J. Heteroscedastic regression in robust econometrics. In Paganoni, A.M., Sangalli, L.M., Secchi, P., Vantini, S. (Eds.): *Proceedings S.Co.2009, Sixth Conference Complex Data Modeling and Computationally Intensive Statistical Methods for Estimation and Prediction*. Politecnico di Milano, Milano, 231-236, 2009.
8. Kalina, J. Asymptotic Durbin-Watson test for robust regression. *Bulletin of the International Statistical Institute* 62, 3406-3409, 2007.
9. Koenker, R. *Quantile regression*. Cambridge: Cambridge University Press, 2005.
10. Plát, P. Modifikace Whiteova testu pro nejmenší vážené čtverce. In Antoch, J., Dohnal, G. (Eds.): *ROBUST 2004, Proceedings of the 13-th Summer school JČMF*. JČMF, Prague, 291-298, 2004. (In Czech.)
11. Víšek, J.Á. Heteroscedasticity resistant robust covariance matrix estimator. *Bulletin of the Czech Econometric Society* 17 (27), 33-49, 2010.
12. Víšek, J.Á. Robustifikace zobecněné metody momentů. In Kupka, K. (Ed.): *Analýza dat 2004/II, Statistické metody pro technologii a výzkum*. Pardubice: Trilobyte Statistical Software, Pardubice, 171-193, 2004. (In Czech.)
13. Víšek, J.Á. Regression with high breakdown point. In Antoch, J., Dohnal, G. (Eds.): *Proceedings of ROBUST 2000, Summer School of JČMF*. Prague: JČMF and Czech Statistical Society, 324-356, 2001.
14. White, H. A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity. *Econometrica* 48 (4), 817-838, 1980.

**Contact**

Jan Kalina

Institute of Computer Science of the Academy of Sciences of the Czech Republic

Pod Vodarenskou vezi 271/2, Praha, Czech republic

kalina@euromise.cz