# Lecture 4: Relaxing the Assumptions of the Linear Model

### R.G. Pierse

## 1 Overview

### 1.1 The Assumptions of the Classical Model

$$\mathrm{E}(u_i) = 0 \quad , \quad i = 1, \cdots, n \qquad (A1)$$

$$\mathrm{E}(u_i^2) = \sigma^2 \quad , \quad i = 1, \cdots, n \qquad (A2)$$

$$\mathrm{E}(u_i u_j) = 0 \quad , \quad i, j = 1, \cdots, n \quad j \neq i \qquad (A3)$$

X values are *fixed in repeated sampling* $\qquad$ (A4)

The variables $X_j$ are *not perfectly collinear.* $\qquad$ (A4$'$)

$u_i$ is distributed with the *normal distribution* $\qquad$ (A5)

In this lecture we look at the implications of relaxing two of these assumptions: A2 and A3. Assumption A2 is the assumption of *homoscedasticity*, that the error variance is constant over all observations. If this assumption does not hold then the errors are *heteroscedastic* and

$$\mathrm{E}(u_i^2) = \sigma_i^2 \quad , \quad i = 1, \cdots, n$$

where the subscript $i$ on $\sigma_i^2$ indicates that the error variance can be different for each observation.

Assumption A3 is the assumption that the errors are *serially uncorrelated*. If this assumption does not hold then we say that the errors are *serially correlated*, or equivalently, that they exhibit *autocorrelation*. Symbolically

$$\mathrm{E}(u_i u_j) \neq 0 \quad j \neq i$$

# 2 Autocorrelation

Here we consider relaxing Assumption A3 and allowing the errors to exhibit *autocorrelation*. Symbolically

$$E(u_i u_j) \neq 0 \quad j \neq i \ . \tag{2.1}$$

Autocorrelation really only makes sense in time-series data where there is a natural ordering for the observations. Hence we assume for the rest of this section that we are dealing with time-series data and use the suffices $t$ and $s$, rather than $i$ and $j$ to denote observations.

## 2.1 The First Order Autoregressive Model

The assumption (2.1) is too general to deal with as it stands and we need to have a more precise model of the form that the autocorrelation takes. Specifically, we consider the hypothesis that the errors follow a *first-order autoregressive* or *AR(1)* scheme

$$u_t = \rho u_{t-1} + \varepsilon_t \quad , \quad t = 1, \cdots, T \tag{2.2}$$

$$-1 < \rho < 1$$

where $u_t$ and $\varepsilon_t$ are assumed to be independent error processes and $\varepsilon_t$ has the standard properties:

$$E(\varepsilon_t) = 0 \quad , \quad E(\varepsilon_t^2) = \sigma_\varepsilon^2 \quad , \quad t = 1, \cdots, T$$

and

$$E(\varepsilon_t \varepsilon_s) = 0 \quad , \quad t, s = 1, \cdots, T \quad s \neq t \ .$$

By successive substitution in (2.2) we can write

$$u_t = \varepsilon_t + \rho \varepsilon_{t-1} + \rho^2 \varepsilon_{t-2} + \cdots + \rho^{t-1} \varepsilon_1 + u_0$$

where the initial value $u_0$ is taken as *fixed* with $u_0 = 0$. Hence it follows that $E(u_t) = 0$.

Consider the variance of $u_t$ where $u_t$ follows a first order autoregressive process:

$$\begin{aligned} E(u_t^2) &= E(\rho u_{t-1} + \varepsilon_t)^2 \\ &= \rho^2 E(u_{t-1}^2) + E(\varepsilon_t^2) + 2\rho E(u_{t-1}\varepsilon_t) \end{aligned}$$

but the final term is zero since $u_t$ and $\varepsilon_t$ are assumed independent and, in the first term, $E(u_{t-1}^2) = E(u_t^2)$ since assumption A2 is still assumed to hold, so that

$$E(u_t^2) = \frac{\sigma_\varepsilon^2}{1 - \rho^2} \ .$$

Similarly, the first order covariance

$$
\begin{aligned}
\mathrm{E}(u_t u_{t-1}) &= \mathrm{E}(\rho u_{t-1} + \varepsilon_t)\, u_{t-1} \\
&= \rho\, \mathrm{E}(u_{t-1}^2) + \rho\, \mathrm{E}(u_{t-1}\varepsilon_t) \\
&= \frac{\sigma_\varepsilon^2 \rho}{1 - \rho^2}
\end{aligned}
$$

and, more generally,

$$
\mathrm{E}(u_t u_{t-s}) = \frac{\sigma_\varepsilon^2 \rho^s}{1 - \rho^2} \quad , \quad s \geq 0 \ .
$$

The autocorrelation between $u_t$ and $u_{t-s}$ decreases as the distance between the observations, $s$, increases. If $\rho$ is negative, then this autocorrelation alternates in sign.

## 2.2 Generalised Least Squares

Consider the multiple regression model

$$
Y_t = \beta_1 + \beta_2 X_{2t} + \cdots + \beta_k X_{kt} + u_t \quad , \quad t = 1, \cdots, T \tag{2.3}
$$

Lagging this equation and multiplying by $\rho$ gives

$$
\rho Y_{t-1} = \rho \beta_1 + \rho \beta_2 X_{2,t-1} + \cdots + \rho \beta_k X_{k,t-1} + \rho u_{t-1} \tag{2.4}
$$

and, subtracting from the original equation,

$$
Y_t - \rho Y_{t-1} = (1-\rho)\beta_1 + \beta_2(X_{2t} - \rho X_{2,t-1}) + \cdots + \beta_k(X_{kt} - \rho X_{k,t-1}) + u_t - \rho u_{t-1}
$$

or

$$
Y_t^* = \beta_1^* + \beta_2 X_{2t}^* + \cdots + \beta_k X_{kt}^* + \varepsilon_t \quad , \quad t = 2, \cdots, T \tag{2.5}
$$

where $Y_t^* = Y_t - \rho Y_{t-1}$ and $X_{jt}^* = X_{jt} - \rho X_{j,t-1}$ are transformed variables known as *quasi-differences*. By transforming the equation, the error process has been tranformed into one that obeys all the classical assumptions so that, if the value of $\rho$ is known, estimation of (2.5) gives the BLUE of the model. This estimator is known as the *Generalised Least Squares* or *GLS* estimator. Note that by quasi-differencing we lose one observation from the sample, because the first observation cannot be quasi-differenced. As long as $T$ is large enough, we don't need to worry about this. For the simple regression model with only one explanatory variable

$$
\widetilde{b} = \frac{\sum_{t=2}^{T}(x_t - \rho x_{t-1})(y_t - \rho y_{t-1})}{\sum_{t=2}^{T}(x_t - \rho x_{t-1})^2}
$$

is the *GLS* estimator and

$$
Var(\widetilde{b}) = \frac{\sigma_\varepsilon^2}{\sum_{t=2}^{T}(x_t - \rho x_{t-1})^2} \ .
$$

## 2.3 Estimating $\rho$: The Cochrane-Orcutt Procedure

In practice, the autoregressive coefficient $\rho$ is unknown and so has to be estimated. The Cochrane-Orcutt (1949) method is an iterative procedure to implement a *feasible GLS* estimator and estimate both $\rho$ and the regression parameters $\beta$ efficiently.

**Step 1**: Estimate the regression equation

$$Y_t = \beta_1 + \beta_2 X_{2t} + \cdots + \beta_k X_{kt} + u_t \quad , \quad t = 1, \cdots, T$$

by OLS and form the OLS residuals $e_t$.

**Step 2:** Run the regression

$$e_t = \rho e_{t-1} + v_t .$$

to give the OLS estimator

$$\widehat{\rho} = \frac{\sum_{t=2}^{T} e_t e_{t-1}}{\sum_{t=2}^{T} e_t^2}$$

**Step 3**: Form the quasi-differenced variables $Y_t^+ = Y_t - \widehat{\rho} Y_{t-1}$, and $X_{j,t}^+ = X_{j,t} - \widehat{\rho} X_{j,t-1}$, $j = 2, \cdots, k$ and run the regression

$$Y_t^+ = \beta_1^+ + \beta_2 X_{2t}^+ + \cdots + \beta_k X_{kt}^+ + \varepsilon_t^+ \quad , \quad t = 2, \cdots, T$$

This gives a new set of estimates for the $\beta$ parameters, and a new set of residuals, $e_t^+$.

Steps 2 and 3 are then executed repeatedly to derive new estimates of $\beta$ and $\rho$ and this process continues until there is no change in the estimate of $\rho$ obtained from successive iterations. This is the full Cochrane-Orcutt iterative procedure. However, a simplified version is the *two-step procedure* which stops after the second step when the OLS estimator of $\rho$ has been obtained.

## 2.4 Properties of OLS

In the presence of autocorrelation, OLS is no longer BLUE. However, it remains unbiased since

$$\mathrm{E}(\widehat{b}) = \mathrm{E}\left(\frac{\sum_{t=1}^{T} x_t y_t}{\sum_{t=1}^{T} x_t^2}\right) = b + \mathrm{E}\left(\frac{\sum_{t=1}^{T} x_t u_t}{\sum_{t=1}^{T} x_t^2}\right) = b$$

where the last equality uses the result that $\mathrm{E}(u_t) = 0$.

Because OLS is no longer BLUE, it follows that the OLS variance must be larger than that of the BLUE *GLS* estimator. Hence, in hypothesis testing, standard errors will be larger than they should be, so that coefficients will appear less significant than they truly are. This will lead to *incorrect inference*.

## 2.5 Testing for Autocorrelation

In practice we do not know whether or not we have autocorrelation. If autocorrelation is ignored, then although parameter estimates remain unbiased, inference based on those estimates will be incorrect. It is thus important to be able to test for autocorrelation in the OLS model. If autocorrelation is detected, then we can re-estimate the model by GLS.

### 2.5.1 The Durbin-Watson test

This is the most famous test for autocorrelation. It is an *exact* test, in that the distribution of the test statistic holds for all sizes of sample. However, it is only valid when (a) the regressors $X$ are fixed in repeated samples and (b) an intercept is included in the regression. Note in particular that the first assumption is violated if the regressors include any lagged dependent variables. The Durbin-Watson (1951) statistic is given by the formula

$$d = \frac{\sum_{t=2}^{T}(e_t - e_{t-1})^2}{\sum_{t=2}^{T} e_t^2} \ . \tag{2.6}$$

Expanding the numerator gives

$$\begin{aligned} d &= \frac{\sum_{t=2}^{T} e_t^2 + \sum_{t=2}^{T} e_{t-1}^2 - 2\sum_{t=2}^{T} e_t e_{t-1}}{\sum_{t=2}^{T} e_t^2} \\ &\simeq 2\left(1 - \frac{\sum_{t=2}^{T} e_t e_{t-1}}{\sum_{t=2}^{T} e_t^2}\right) = 2(1 - \widehat{\rho}) \end{aligned}$$

since $\sum_{t=2}^{T} e_{t-1}^2 \equiv \sum_{t=1}^{T-1} e_t^2 \simeq \sum_{t=2}^{T} e_t^2$. Note that $0 \leq d \leq 4$, with $\mathrm{E}(d) = 2$, and $d < 2$ indicative of positive autocorrelation, and $d > 2$ indicative of negative autocorrelation.

In general, the distribution of $d$ is a function of the explanatory variables $X$, so to compute the exact distribution of the statistic is very complicated. However, Durbin and Watson were able to show that the distribution of $d$ is bounded by the distributions of two other statistics $d_L$ and $d_U$ which *do not depend on $X$*. Critical values of the distributions of $d_L$ and $d_U$ were tabulated by Durbin and Watson.

The procedure for applying the Durbin-Watson test for testing positive autocorrelation is as follows:

$$\begin{aligned} d &< d_L : \text{reject the null hypothesis of no positive autocorrelation} \\ d_L &\leq d \leq d_U : \text{inconclusive region} \\ d &> d_U : \text{do not reject the null hypothesis of no autocorrelation} \end{aligned}$$

Similarly, to test for negative autocorrelation:

$$d > 4 - d_L : \text{reject the null hypothesis of no negative autocorrelation}$$
$$4 - d_U \leq d \leq 4 - d_L : \text{inconclusive region}$$
$$d < 4 - d_U : \text{do not reject the null hypothesis of no autocorrelation}$$

In both cases, there is an inconclusive region where the test does not allow us to say whether or not there is autocorrelation. For example, with $T = 25$ and $k = 2$, the 5% critical values are $d_L = 1.29$ and $d_U = 1.45$ so that, for any value of $d$ in this range, the result of the test is indeterminate.

The inconclusive region in the Durbin-Watson test is clearly a nuisance. One solution is to adopt the *modified d test*, where test values in the inconclusive region are treated as rejections of the null hypothesis. The justification for this is that, in general, with the types of variable encountered in economics, $d_U$ is likely to be closer to the true distributionof $d$ than $d_L$. It can also be regarded as adopting a cautious strategy with respect to autocorrelation.

### 2.5.2 Testing for Higher Order AR Processes: An LM test

Suppose that the error term $u_t$ is generated by the $p$th order autoregressive process

$$u_t = \rho_1 u_{t-1} + \rho_2 u_{t-2} + \cdots + \rho_p u_{t-p} + \varepsilon_t \quad , \quad t = 1, \cdots, T$$

A test of the null hypothesis:

$$H_0 : \rho_1 = \rho_2 = \cdots = \rho_p = 0$$

can be constructed by the following procedure:

    1.     Estimate the regression model by OLS and obtain the residuals $e_t$.

    2.     Regress the residuals $e_t$ on all the regressors, plus the lagged residuals, $e_{t-1}, e_{t-2}, \cdots, e_{t-p}$.

    3.     Obtain the $R^2$ from this auxiliary regression.

Then it can be shown that

$$(T - p)R^2 \sim_a \chi_p^2$$

is a test of $H_0$ that is valid *asymptotically*, i.e. as $T \to \infty$. In practice, this test will be approximately valid as long as the sample size is 'large'. Note that, in order to perform the auxiliary regression in step 2, the first $p$ observations need to be dropped.

This test was derived independently by Breusch (1978) and Godfrey (1978) and is sometimes called the Breusch-Godfrey or the *Lagrange Multiplier* ($LM$) test for $p$th order autocorrelation. Note that this test does not depend on the assumption of fixed regressors and so is still valid when the regressors include lagged dependent variables.

## 2.6  Other Models of Autocorrelation

So far we have considered only autoregressive models of autocorrelation. One other important model of autocorrelation is the $q$th order *moving average* model $MA(q)$:

$$u_t = \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \cdots + \theta_q \varepsilon_{t-q} \quad , \quad t = 1, \cdots, T$$

This model has the property that

$$\begin{aligned}
\mathrm{E}(u_t u_{t-s}) &= \sigma_\varepsilon^2 \sum_{j=s}^{q} \theta_j \theta_{j-s} \quad , \quad 0 \leq s \leq q \\
&= 0 \quad , \quad s > q
\end{aligned}$$

so that there is no autocorrelation after the $q$th order. $MA$ models of autocorrelation are generally more difficult to estimate than $AR$ models. It can be shown that the Breusch-Godfrey Lagrange Multiplier test of the previous section is also a valid test of the hypothesis of no autocorrelation in the $MA(p)$ model.

Finally, one other model of autocorrelation is the mixed autoregressive-moving average model, or $ARMA(p, q)$ model. This is a very general model of autocorrelation.

# 3  Heteroscedasticity

Where autocorrelation is a problem of time-series data, heteroscedasticity is primarily a problem of cross-sectional data. Here we consider relaxing assumption A2 to allow the error variance to differ between observations, or symbolically,

$$\mathrm{E}(u_i^2) = \sigma_i^2 = k_i^2 \sigma^2 \quad , \quad i = 1, \cdots, n$$

where $k_i$ are fixed constants.

## 3.1  Generalised Least Squares

Consider again, the multiple regression model

$$Y_i = \beta_1 + \beta_2 X_{2i} + \cdots + \beta_k X_{ki} + u_i \quad , \quad i = 1, \cdots, n \tag{3.1}$$

Dividing through by $k_i$ gives

$$\frac{Y_i}{k_i} = \frac{\beta_1}{k_i} + \beta_2 \frac{X_{2i}}{k_i} + \cdots + \beta_k \frac{X_{ki}}{k_i} + \frac{u_i}{k_i} \tag{3.2}$$

or

$$Y_i^* = \beta_1^* + \beta_2 X_{2i}^* + \cdots + \beta_k X_{ki}^* + u_i^* \quad , \quad i = 1, \cdots, n \tag{3.3}$$

where $Y_i^* = Y_i/k_i$ and $X_{ji}^* = X_{ji}/k_i$ are transformed variables and

$$Var(u_i^*) = \text{E}\left(\frac{u_i}{k_i}\right)^2 = \frac{1}{k_i^2}\text{E}(u_i^2) = \sigma^2$$

so that the transformed error process is homoscedastic and obeys all the classical assumptions. As long as the weights $k_i$ are known, estimation of (3.3) gives the BLUE estimator of the heteroscedastic model. This estimator is the *Generalised Least Squares* or *GLS* estimator. It is also known as the *Weighted Least Squares* or *WLS* estimator.

## 3.2 The Properties of OLS

In the presence of heteroscedasticty, OLS is no longer BLUE. However, it is still unbiased since $\text{E}(u_i) = 0$. As with autocorrelation, the effect of heteroscedasticity is that the estimated variance is larger than that of the 'correct' model so that inference based on OLS estimates will be incorrect.

## 3.3 Testing for Heteroscedasticity

There are several tests of heteroscedasticity available, based on different assumptions about the form that the heteroscedasticity takes.

### 3.3.1 The Goldfeld-Quandt Test

Suppose that the heteroscedasticity is proportional to the square of one of the regressors $X_j$:

$$\sigma_i^2 = \sigma^2 X_{j,i}^2$$

The Goldfeld-Quandt (1972) test is based on the following procedure:

    1.    Order the observations by the variable $X_j$.

    2.    Omit the central $c$ observations and divide the sample into two groups, each of $(n-c)/2$ observations.

    3.    Fit separate regressions to the two sub-samples, obtaining the Residual Sum of Squares $RSS_1$ and $RSS_2$ respectively, where $RSS_2$ is the sub-sample corresponding to the largest observations of $X_j$.

    Then on the null hypothesis of homoscedasticity

$$\frac{RSS_2}{RSS_1} \sim F_{(n-c-2k)/2,(n-c-2k)/2}$$

This test is an *exact* test. When there is more than one regressor, then the choice of the variable on which to order: $X_j$ is arbitrary, and the assumption that heteroscedasticity is related to a single regressor is less appealing.

8

### 3.3.2 The Breusch-Pagan-Godfrey Test

This test is based on the assumption that the heteroscedasticity takes the form

$$\sigma_i^2 = \sigma^2 + \alpha_1 Z_{1i} + \cdots + \alpha_m Z_{mi}$$

where the $Z_j$'s are non-stochastic variables. Some or all of the $X_j$'s may be included in $Z$. The null hypothesis of homoscedasticity is equivalent to a test of the hypothesis that

$$H_0 : \alpha_1 = \cdots = \alpha_m = 0 \ .$$

The test is based on the following procedure:
  1    Estimate the regression model by OLS and obtain the residuals $e_i$.
  2.    Construct the variables $p_i = e_i^2 / \tilde{\sigma}^2$ where $\tilde{\sigma}^2 = \sum_{i=1}^n e_i^2 / n$.
  3.    Run the regression

$$p_i = \alpha_0 + \alpha_1 Z_{1i} + \cdots + \alpha_m Z_{mi} + v_i$$

  and obtain the explained sum of squares $ESS$.
  Then, it can be shown that

$$\frac{ESS}{2} \sim_a \chi_m^2$$

is a test of $H_0$ that is valid *asymptotically*, i.e. as $n \to \infty$. In practice, this test will be approximately valid as long as the sample size is 'large'.

### 3.3.3 The White test for general heteroscedasticity

This is a test for general heteroscedasticity of unknown form developed by White (1980). It is based on a regression of the squared residuals $e_t^2$ on all (non-redundant) cross-products of the regressors $X_{ji}X_{li}$, for all $j = 1, \cdots, k$, and $l = 1, \cdots, k$. The auxiliary regression takes the form:

$$\begin{aligned} e_i^2 &= \gamma_1 X_{1i}^2 + \gamma_2 X_{1i}X_{2i} + \cdots + \gamma_k X_{1i}X_{ki} \\ &\quad + \gamma_{k+1}X_{2i}^2 + \gamma_{k+2}X_{2i}X_{3i} + \cdots + \gamma_{2k-1}X_{2i}X_{ki} + \cdots + \gamma_{k(k+1)/2}X_{ki}^2 + u_i. \end{aligned}$$

Note that since $X_{1i}$ is an intercept, $X_{1i}^2 = X_{1i}$ and the first $k$ terms are just the original equation regressors. The test statistic is based on the $R^2$ from this auxiliary regression and is asymptotically valid. White (1980) showed that, on the null hypothesis of no heteroscedasticity,

$$nR^2 \sim_a \chi_{k(k+1)/2-1}^2.$$

When there are a large number of regressors, the number of regressor cross-products $k(k+1)/2$ can quickly become very large and this test may be inpracticable. An alternative version of the test using only the squared regressor terms $X_{ji}^2$ is also available and takes the form

$$e_i^2 = \gamma_1 X_{1i}^2 + \gamma_2 X_{2i}^2 + \cdots + \gamma_k X_{ki}^2 + u_i.$$

For this case we can show that, on the null hypothesis of no heteroscedasticity,

$$nR^2 \sim_a \chi_{k-1}^2.$$

## 3.4 White heteroscedasticity-consistent standard errors

The problem with OLS estimation under heteroscedasticity is that the estimated standard errors are incorrect so that inference is invalid. Efficient estimation requires knowing the form of the heteroscedasticity. White (1980) derived an estimator for the variance of the OLS coefficients which remains consistent under the hypothesis of general heteroscedasticity. It does not require that the form of the heteroscedasticity be known. White standard errors are computed by most regression packages. The formula, for the simple regression case, is

$$\widetilde{Var}(\widehat{b}) = \frac{\sum_{i=1}^{n} x_i^2 e_i^2}{\left(\sum_{i=1}^{n} x_i^2\right)^2}$$

## 3.5 ARCH: Heteroscedasticity in Time Series Models

As a rule, heteroscedasticity is a cross-sectional, rather than a time-series problem. However, Engle (1982) has proposed the following model for heteroscedasticity in a time-series model:

$$Var(u_t) = \sigma_t^2 = \alpha_0 + \alpha_1 u_{t-1}^2 + \cdots + \alpha_p u_{t-p}^2$$

This is known as the *pth order Autoregressive Conditional Heteroscedastic* model or $ARCH(p)$ model. A test for homoscedasticity in this model is based on the $R^2$ from the following regression of the squared OLS residuals

$$\widehat{u}_t^2 = \widehat{\alpha}_0 + \widehat{\alpha}_1 \widehat{u}_{t-1}^2 + \cdots + \widehat{\alpha}_p \widehat{u}_{t-p}^2 .$$

It can be shown that

$$TR^2 \sim_a \chi_p^2$$

This is another example of a Lagrange Multiplier test whose distribution is valid asymptotically.

# 4 Autocorrelation as a Symptom of Misspecification

Both autocorrelation and heteroscedasticity lead to a systematic pattern in the OLS residuals. The tests developed above are based on looking for particular systematic patterns in the residuals as evidence for autocorrelation or for heteroscedasticity. However, finding such a pattern may be evidence of something else. In particular, omitting a variable from a regression, the most common form of *misspecification*, will often lead to autocorrelation in the equation residuals, if the omitted variable is itself autocorrelated, and this is generally the case in economic variables. This suggests that tests for autocorrelation can also be interpreted as more general tests of misspecification, and that finding evidence of serial correlation should not necessarily lead to the automatic adoption of the GLS estimator but to a more general reconsideration of the specification of the equation. An apposite quotation reflects the current view among econometricians:

> There is no universally effective way of avoiding misinterpreting misspecification of the regression function as the presence of serially correlated errors. R. Davidson and J.G. MacKinnon (1993)

# References

[1] Breusch, T.S. (1978),'Testing for autocorrelation in dynamic linear models', *Australian Economic Papers*, 17, 334–355.

[2] Breusch, T.S. and A. Pagan (1979), 'A simple test for heteroscedasticity and random coefficient variation', *Econometrica*, 47, 1287–1294.

[3] Cochrane, D. and G.H. Orcutt, (1949), 'Application of least squares regressions to relationships containing autocorrelated error terms', *Journal of the American Statistical Association*, 44, 32–61.

[4] Davidson, R. and J.G. MacKinnon (1993), *Estimation and Inference in Econometrics*, Oxford University Press, New York.

[5] Durbin, J. and G.S. Watson (1951), 'Testing for serial correlation in least squares regression', *Biometrika*, 38, 159–171.

[6] Engle, R.F. (1982), 'Autoregressive conditional heteroscedasticity with estimates of the variance of United Kingdom inflation', *Econometrica*, 50, 987–1007.

[7] Godfrey, L.G. (1978), 'Testing against general autoregressive and moving average error models when the regressors include lagged dependent variables', *Econometrica*, 46, 1293–1302.

[8] Godfrey, L.G. (1978), 'Testing for multiplicative heteroscedasticity', *Journal of Econometrics*, 8, 227–236.

[9] Goldfeld, S.M. and R.E. Quandt (1972), *Nonlinear Methods in Econometrics*, North-Holland, Amsterdam.

[10] White, H. (1980), 'A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity', *Econometrica*, 48, 817–838.