

# IIA-3 Econometrics: Supervision 4

Emre Usenmez

Christmas Break 2024

## Topics Covered

**Faculty Qs:** Endogeneity; proxy variable; instrumental variable

**Supplementary Qs:** Heteroskedasticity; White's robust standard errors; Goldfeld-Quandt test; Breusch-Pagan test; Lagrange Multiplier test; White test; autocorrelation; Durbin-Watson  $d$  test; Breusch-Godfrey test; Durbin's alternative test; Prais-Winsten transformation; Cochrane-Orcutt Two-Step Procedure

## Related Reading:

Dougherty, *Introduction to Econometrics*, 5<sup>th</sup> ed, OUP

Chapter 7: Heteroskedasticity

Chapter 8: Stochastic Regressors and Measurement Errors

Chapter 9: Simultaneous Equations Estimation

Chapter 12: Autocorrelation

Wooldridge J M (2021) *Introductory Econometrics: A Modern Approach*, 7<sup>th</sup> ed,

Chapter 8: Heteroskedasticity

Chapter 9: More on Specification and Data Issues

Chapter 12: Serial Correlation and Heteroskedasticity in Time Series Regressions

Chapter 16: Simultaneous Equations Models

Gujarati, D N and Porter, D (2009) *Basic Econometrics* 7<sup>th</sup> International ed, McGraw Hill

Chapter 11 Heteroscedasticity: What Happens If the Error Variance is Nonconstant?

Chapter 12: Autocorrelation: What Happens If the Error Terms are Correlated?

Chapter 13: Econometric Modeling: Model Specification and Diagnostic Testing

Chapter 20: Simultaneous-Equation Methods

Very grateful to Dr Oleg Kitov and Dr Clive Lawson for the very informative stylized answers to previous iterations of the supervision questions.

## FACULTY QUESTIONS

### QUESTION 1

Consider the following bivariate linear regression

$$y = \alpha + T\beta + u$$

where  $T$  is a binary treatment regressor,  $\alpha$  and  $\beta$  are unknown parameters, and  $u$  is an error term.

(a) Describe in two sentences an empirical, real-life example where such an equation might arise.

**Answer:** We can think of  $T$  as "graduated from university" and  $y$  as "annual earning after 10 years of graduation."

---

(b) Why might  $u$  be heteroskedastic in your example.

**Answer:** The variance of earnings will likely to be smaller across people who did not graduate from a university compared to those who did it. This may be because those who did not go to university are less likely to be in the professions such as lawyers or doctors, and more likely to be in lower-paying jobs, or unemployed, or out of labor force.

---

(c) Why might  $T$  be endogenous in your example?

**Answer:** Broadly, variables that are correlated with the error term are called *endogeneous variables*, and those that are uncorrelated with the error term are called *exogeneous variables*.<sup>1</sup>

Thus the question is asking us to consider some of the reasons as to why  $T$  might be correlated with the error term. There are certainly nonnegligible number of high earners who either never went to a university or dropped out. There may be omitted variable or even simultaneity is possible.

Let's consider what the implications of of endogeneity are for the OLS estimator of  $\beta$ .

Variable  $T$  would be endogenous if  $\mathbb{E}(u|T) \neq 0$ . Endogeneity would imply that  $Cov(T, u) \neq 0$ .

We can first look at whether it is biased. For that, we need to use the law of iterated expectations whereby

$$\mathbb{E}(\hat{\beta}^{OLS}) = \mathbb{E}[\mathbb{E}(\hat{\beta}^{OLS} | T_1, \dots, T_n)]$$

The OLS estimator of  $\beta$  would be:

$$\begin{aligned} \mathbb{E}(\hat{\beta}^{OLS} | T_1, \dots, T_n) &= \mathbb{E}\left(\frac{\widehat{Cov}(T_i, Y_i)}{\widehat{Var}(T_i)} \mid T_1, \dots, T_n\right) = \mathbb{E}\left(\frac{\hat{\sigma}_{TY}}{\hat{\sigma}_{TT}} \mid T_1, \dots, T_n\right) \\ &= \mathbb{E}\left(\frac{\sum_{i=1}^n (T_i - \bar{T})(Y_i - \bar{Y})}{\sum_{i=1}^n (T_i - \bar{T})^2} \mid T_1, \dots, T_n\right) \\ &= \mathbb{E}\left(\frac{\sum_{i=1}^n (T_i - \bar{T})((\alpha + \beta T_i + u_i) - (\alpha + \beta \bar{T} + \bar{u}))}{\sum_{i=1}^n (T_i - \bar{T})^2} \mid T_1, \dots, T_n\right) \\ &= \mathbb{E}\left(\frac{\sum_{i=1}^n (T_i - \bar{T})(\beta(T_i - \bar{T}) + u_i - \bar{u})}{\sum_{i=1}^n (T_i - \bar{T})^2} \mid T_1, \dots, T_n\right) \\ &= \mathbb{E}\left(\frac{\sum_{i=1}^n \beta(T_i - \bar{T})^2 + \sum_{i=1}^n (T_i - \bar{T})(u_i - \bar{u})}{\sum_{i=1}^n (T_i - \bar{T})^2} \mid T_1, \dots, T_n\right) \\ &= \mathbb{E}\left(\beta + \frac{\sum_{i=1}^n (T_i - \bar{T})(u_i - \bar{u})}{\sum_{i=1}^n (T_i - \bar{T})^2} \mid T_1, \dots, T_n\right) \end{aligned}$$

<sup>1</sup>See Chapter 12: Instrumental Variables Regression p.428 in Stock J H, and Watson M W (2020) Introduction to Econometrics, 4<sup>th</sup> Global Ed, Pearson; and Section 8.5: Instrumental Variables in Dougherty C (2016) Introduction to Econometrics 5<sup>th</sup> ed, OUP in addition to Chapter 9: More on Specification and Data Issues in Wooldridge J M (2021) Introductory Econometrics: A Modern Approach, 7<sup>th</sup> ed, Cengage

$$\begin{aligned}
&= \mathbb{E} \left( \beta + \frac{\sum_{i=1}^n (T_i - \bar{T}) u_i - \bar{u} \sum_{i=1}^n (T_i - \bar{T})}{\sum_{i=1}^n (T_i - \bar{T})^2} \middle| T_1, \dots, T_n \right) \\
&= \mathbb{E} \left( \beta + \frac{\sum_{i=1}^n (T_i - \bar{T}) u_i - \bar{u} \left( \sum_{i=1}^n T_i - n\bar{T} \right)}{\sum_{i=1}^n (T_i - \bar{T})^2} \middle| T_1, \dots, T_n \right) \\
&= \mathbb{E} \left( \beta + \frac{\sum_{i=1}^n (T_i - \bar{T}) u_i - \bar{u} (n\bar{T} - n\bar{T})}{\sum_{i=1}^n (T_i - \bar{T})^2} \middle| T_1, \dots, T_n \right) \\
&= \mathbb{E} \left( \beta + \frac{\sum_{i=1}^n (T_i - \bar{T}) u_i}{\sum_{i=1}^n (T_i - \bar{T})^2} \middle| T_1, \dots, T_n \right) \\
&= \beta + \frac{\sum_{i=1}^n (T_i - \bar{T}) \mathbb{E}(u_i \mid T_1, \dots, T_n)}{\mathbb{E} \left( \sum_{i=1}^n (T_i - \bar{T})^2 \mid T_1, \dots, T_n \right)}
\end{aligned}$$

Notice that since  $\mathbb{E}(u|T) \neq 0$ , the numerator of this last expression is also nonzero. That is,  $\sum_{i=1}^n (T_i - \bar{T}) \mathbb{E}(u_i \mid T_1, \dots, T_n) \neq 0$ . Therefore the expectation of this expectation is also not equal to  $\beta$ :

$$\mathbb{E}(\hat{\beta}^{OLS}) = \mathbb{E} \left[ \mathbb{E}(\hat{\beta}^{OLS} \mid T_1, \dots, T_n) \right] = \mathbb{E} \left[ \mathbb{E} \left( \beta + \frac{\sum_{i=1}^n (T_i - \bar{T}) u_i}{\sum_{i=1}^n (T_i - \bar{T})^2} \middle| T_1, \dots, T_n \right) \right] \neq \beta$$

which means the OLS estimator is *not* unbiased.

We can also check for consistency by examining the probability limit of this expression as  $n$  tends towards infinity. For that, we can rewrite the OLS estimator as:

$$\hat{\beta}^{OLS} = \beta + \frac{\frac{1}{n} \sum_{i=1}^n (T_i - \bar{T}) u_i}{\frac{1}{n} \sum_{i=1}^n (T_i - \bar{T})^2}$$

Using the law of large numbers, we can see that as  $n \rightarrow \infty$

$$\frac{1}{n} \sum_{i=1}^n (T_i - \bar{T}) u_i \xrightarrow{p} \mathbb{E}[(T_i - \bar{T}) u_i] = \text{Cov}(T_i, u_i) \neq 0$$

and

$$\frac{1}{n} \sum_{i=1}^n (T_i - \bar{T})^2 \xrightarrow{p} \mathbb{E}[(T_i - \bar{T})^2] = \text{Var}(T_i) = \sigma_T^2 < \infty$$

Note that  $\text{Var}(T_i) = \sigma_T^2 < \infty$  is an additional assumption.

Since  $\text{Cov}(T_i, u_i) \neq 0$ , the OLS estimator as  $n \rightarrow \infty$  (using Slutsky's theorem):

$$\hat{\beta}^{OLS} \xrightarrow{p} \beta + \frac{\text{Cov}(T_i, u_i)}{\text{Var}(T_i)} \neq \beta$$

which means the OLS estimator is not only biased but also inconsistent for  $\beta$ .

(d) Suppose a single instrument  $z$  is available. Show that the population coefficient  $\beta$  satisfies

$$\beta = \frac{\text{Cov}(z, y)}{\text{Cov}(z, T)}$$

where  $\text{Cov}(z, y)$  and  $\text{Cov}(z, T)$  denote, respectively, the population covariance between  $z$  and  $y$ , and  $z$  and  $T$ . How can you use this information to obtain a consistent estimate of  $\beta$ ?

**Answer:** Instrument  $z$  needs to satisfy the following conditions:

- *Instrument relevance:*  $z$  must have non-trivial explanatory power for  $T$ , namely  $\text{Cov}(z, T) \neq 0$ .
- *Instrument exogeneity:*  $z$  must affect  $Y$  only through its influence on  $T$  and not in any other way. That is,  $z$  must be exogenous with respect to  $u$  in regression  $y = \alpha + \beta T + u$ . Formally,  $\mathbb{E}(u|z) = 0$ . This is why it is said " $z$  is exogenous in  $y = \alpha + \beta T + u$ ". Exogeneity of instrument  $z$  implies that  $\text{Cov}(z, u) = 0$ .

In the context of omitted variables, instrument exogeneity means that  $z$  should be uncorrelated with the omitted variables, i.e.  $\text{Cov}(z, u) = 0$ , and  $z$  should be related, positively or negatively, to the endogenous explanatory variable  $T$ , i.e.  $\text{Cov}(z, T) \neq 0$ .<sup>2</sup>

The underlying reasoning is that if an instrument is relevant, then variation in that instrument  $z$  is related to variation in  $T$ , and if it is also exogenous, then that part of the variation of  $T$  captured by  $z$  is exogenous. Therefore, an instrument that is relevant and exogenous can capture movements in  $T$  that are exogenous. This exogenous variation can in turn be used to estimate the population coefficient  $\beta$ .<sup>3</sup>

These conditions serve to *identify* the parameter  $\beta$ . In this context, *identification of a parameter* means that we can write  $\beta$  in terms of population moments that can be estimated using a sample of data.

To write  $\beta$  in terms of population covariances we use  $y = \alpha + \beta T + u$ :

$$\text{Cov}(z, y) = \text{Cov}(z, \alpha + \beta T + u) = \beta \text{Cov}(z, T) + \text{Cov}(z, u)$$

<sup>2</sup>see Section 15-1: Omitted Variables in a Simple Regression Model in Wooldridge J M (2021) Introductory Econometrics: A Modern Approach, 7<sup>th</sup> ed, Cengage

<sup>3</sup>see Section 12.1: The IV Estimator with a Single Regressor and a Single Instrument in Stock and Watson (2020, 4<sup>th</sup> ed.).

Since instrument exogeneity condition assumes that  $Cov(z, u) = 0$  then  $Cov(z, y) = \beta Cov(z, T)$ . Rearranging this gives:

$$\beta = \frac{Cov(z, y)}{Cov(z, T)}$$

as desired. Notice that this only holds if instrument relevance also holds, since this expression would fail if  $Cov(z, T) = 0$ . What this expression is telling us is that  $\beta$  is identified by the ratio of population covariance between  $z$  and  $y$  to population covariance between  $z$  and  $T$ .

Given a random sample, we estimate the population quantities by the sample analogs:

$$\hat{\beta}^{IV} = \frac{\frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y})}{\frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})(T_i - \bar{T})} = \frac{\sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y})}{\sum_{i=1}^n (z_i - \bar{z})(T_i - \bar{T})}.$$

With a sample data on  $T$ ,  $y$ , and  $z$  we can obtain the IV estimator above. The IV estimator for the intercept  $\alpha$  is  $\alpha = \bar{y} - \hat{\beta}^{IV} \bar{T}$ . Also notice that when  $z = T$ , we get the OLS estimator of  $\beta$ . That is, when  $T$  is exogeneous, it can be used as its own IV, and the IV estimator is then identical to the OLS estimator.

A similar set of steps we used in part (c) will show that IV estimator is consistent for  $\beta$ . That is,  $\text{plim}_{n \rightarrow \infty}(\hat{\beta}) = \beta$ .

Note that, an important feature of IV estimator is that when  $T$  and  $u$  are in fact correlated, and thus instrumental variables estimation is actually needed, it is essentially never unbiased. This means, in small samples, the IV estimator can have a substantial bias, which is one reason why large samples are preferred.

(e) Can you give an example of an instrument in your example? Argue why it might be a sensible IV.

**Answer:** Distance from nearest college can be an example of an instrument, where  $z = 1$  if individual lived near college and 0 otherwise. This may be violated for a number of reasons, though; for e.g. if wealthy parents choose to live near college. This would mean that  $z$  is correlated with unobserved factors that also affect wage, our  $y$ . For any example, exogeneity and relevance conditions need to be checked.

## QUESTION 2

Consider the following wage equation that explicitly recognizes that ability affects  $\log(wage)$

$$\log(wage) = \alpha + \beta_1 educ + \beta_2 ability + u$$

The above model shows explicitly that we would like to hold ability fixed when measuring the returns on education. Assuming that the primary interest is in obtaining a reliable estimate of the slope parameters  $\beta_1$ , and that there is no direct measurement for ability, explain how you would do this using a method based upon a proxy variable and an IV estimator. In doing so evaluate the following statement:

*“whilst IQ is a good candidate as a proxy for variable for ability, it is not a good instrumental variable for educ.”*

**Answer:** This question is essentially aiming to ensure the students understand the difference between proxy variable and instrumental variable.

proxy variable refers to an *observed* variable that is correlated with but not identical to the *unobserved* variable.

instrumental variable refers to a variable that does not appear in the regression, uncorrelated with the error in the equation, and partially correlated with the endogenous explanatory variable in an equation where such endogenous explanatory variable exists.

### Proxy Variable:

Notice in this question *educ* is observed but *ability* is unobserved, and we would not even know how to interpret it's coefficient  $\beta_2$  since 'ability' itself is a vague concept. We can instead use intelligence quotient, or *IQ*, as a proxy for ability as long as *IQ* is correlated with ability. This is captured by the following simple regression:

$$ability = \delta_0 + \delta_2 IQ + v_2$$

where  $v_2$  is an error due to the fact that *ability* and *IQ* are not exactly related. The parameter  $\delta_2$  measures the relationship between *ability* and *IQ*. If  $\delta_2 = 0$  then *IQ* is not a suitable proxy for *ability*.

Note that the intercept  $\delta_0$  allows *ability* and *IQ* to be measured on different scales and thus can be positive or negative. That is, the unobserved *ability* is not required to have the same average value as *IQ* in the population.

In order to use *IQ* to get unbiased, or at least consistent, estimators for  $\beta_1$ , which is the coefficient of *educ*, we would regress  $\log(wage)$  on *educ* and *IQ*. This is called *the plug-in solution to the omitted variables problem* since we plug-in *IQ* for *ability* before running the OLS. However, since *IQ* and *educ* are not the same, we need to check if this does give consistent estimator for  $\beta_1$ .

For the plug-in solution to provide consistent estimator for  $\beta_1$  the following two assumptions need to hold true:

- The error  $u$  is uncorrelated with *educ* and *ability* as well as *IQ*. That is,  $\mathbb{E}(u | educ, ability, IQ) = 0$ . What this means is that *IQ* is irrelevant in the population model which is true by definition since *IQ* is a proxy for *ability*, it is *ability* that directly affects  $\log(wage)$  not *IQ*.
- The error  $v_2$  is uncorrelated with *educ* and *IQ*. For  $v_2$  to be uncorrelated with *educ*, *IQ* needs to be a 'good' proxy for *ability*.

What is meant by 'good' proxy in this sense is that

$$\mathbb{E}(ability | educ, IQ) = \mathbb{E}(ability | IQ) = \delta_0 + \delta_2 IQ.$$

Here, the first equality, which is the most important one, says that once *IQ* is controlled for, the expected value of *ability* does not depend on *educ*. In other words, *ability* has zero correlation with

*educ* once *IQ* is partialled out. Thus the average level of *ability* only changes with *IQ* and not with *educ*.

To see why these two assumptions are enough for the plug-in solution to work, we can rewrite the  $\log(\text{wage})$  equation in the question as:

$$\begin{aligned}\log(\text{wage}) &= \alpha + \beta_1 \text{educ} + \beta_2 \text{ability} + u \\ &= \alpha + \beta_1 \text{educ} + \beta_2(\delta_0 + \delta_2 \text{IQ} + v_2) + u \\ &= (\alpha + \beta_2 \delta_0) + \beta_1 \text{educ} + \beta_2 \delta_2 \text{IQ} + u + \beta_2 v_2 \\ &= (\alpha + \beta_2 \delta_0) + \beta_1 \text{educ} + \beta_2 \delta_2 \text{IQ} + \epsilon \\ &= \gamma_0 + \beta_1 \text{educ} + \gamma_2 \text{IQ} + \epsilon.\end{aligned}$$

Notice that the composite error  $\epsilon$  depends on both the error in the model of interest in the question,  $u$ , and on the error in the proxy variable equation,  $v_2$ . Since both  $u$  and  $v_2$  have zero mean and each is uncorrelated with *educ* and *IQ*,  $\epsilon$  also has zero mean and is uncorrelated with *educ* and *IQ*.

So when we regress  $\log(\text{wage})$  on *educ* and *IQ*, we will not get unbiased estimators of  $\alpha$  and  $\beta_2$ . Instead, we will get unbiased, or at least consistent, estimators of  $\gamma_0, \beta_1$ , and  $\gamma_2$ . The important thing is that we get good estimators of  $\beta_1$ .

In most cases, the estimate of  $\gamma_2$  is more interesting than an estimate of  $\beta_2$  anyway, since  $\gamma_2$  measures the return to  $\log(\text{wage})$  given one more point on *IQ* score.

#### Bias and Multicollinearity when using a proxy

##### When using a proxy variable can still lead to bias?

If the two assumptions above are not satisfied, then using a proxy variable can lead to a bias. To see this, suppose now that *ability* is not only related to *IQ* but also to *educ*:

$$\text{ability} = \delta_0 + \delta_1 \text{educ} + \delta_2 \text{IQ} + v_3$$

where the error  $v_3$  has a zero mean and uncorrelated with *educ* and *IQ*. In the proxy variable discussion above, it was essentially assumed that  $\delta_1 = 0$ . We can re-write  $\log(\text{wage})$  with this plug-in solution:

$$\begin{aligned}\log(\text{wage}) &= \alpha + \beta_1 \text{educ} + \beta_2 \text{ability} + u \\ &= \alpha + \beta_1 \text{educ} + \beta_2(\delta_0 + \delta_1 \text{educ} + \delta_2 \text{IQ} + v_3) + u \\ &= (\alpha + \beta_2 \delta_0) + (\beta_1 + \beta_2 \delta_1) \text{educ} + \beta_2 \delta_2 \text{IQ} + u + \beta_2 v_3\end{aligned}$$

Since the error term  $u + \beta_2 v_3$  has zero mean and is uncorrelated with *educ* and *IQ*, we have  $\text{plim}(\hat{\beta}_1) = \beta_1 + \beta_2 \delta_1$ . If *educ* is partially and positively correlated with *ability*, i.e.  $\delta_1 > 0$ , and if the coefficient of *ability* is positively correlated with  $\log(\text{wage})$ , i.e.  $\beta_2 > 0$ , then  $\text{plim}(\hat{\beta}_1) = \beta_1 + \beta_2 \delta_1 > \beta_1$  giving us an upward bias. That is, in this case where *IQ* is not a good proxy for *ability* but we still use it, then we'd still be getting an upward bias for the coefficient of *educ*. Having said that, the bias is likely to be smaller than if we ignored the problem of omitted ability entirely.

##### What about multicollinearity?

Even if *IQ* is a good proxy for *ability*, using it in a regression that includes *educ* can exacerbate the multicollinearity problem, which, in turn, is likely to lead a less precise estimate of the coefficient for *educ*, i.e.  $\beta_1$ .

However, notice that

- inclusion of *IQ* in the regression means that the part of *ability* explained by *IQ* is removed from the error term, reducing the error variance. This is likely to be reflected in a smaller



standard error of the regression, though that reduction may not happen because of degrees of freedom adjustment.

- if we want to have a less bias for  $\beta_1$ , ie, the estimator of the coefficient for *educ*, then we have to live with increased multicollinearity. This is an important point. Since *educ* and *ability* are thought to be correlated, and if we could include *ability* in the regression, then there would be inevitable multicollinearity caused by the correlation between these two variables. Since *IQ* is a proxy for *ability*, *educ* and *IQ* are also correlated, and a similar argument ensues.

### Instrumental Variable

Suppose now that the proxy variable does not have the required properties for a consistent estimator of  $\beta_1$ . Then we put *ability* in the error term since it is unobserved and we don't have a proxy for it. This leaves us with:

$$\log(wage) = \beta_0 + \beta_1 educ + \epsilon$$

where  $\epsilon$  contains *ability*. If *ability* and *educ* are correlated, then we have a biased and inconsistent estimate of  $\beta_1$ .

However, we can still use this equation as the basis for estimation as long as we can find an instrumental variable for *educ*. For this we can introduce an *instrumental variable*  $z$  which satisfies the "instrument relevance", i.e.  $Cov(z, educ) \neq 0$ , and "instrument exogeneity", i.e.  $Cov(z, \epsilon) = 0$  conditions as discussed in Question 1(d).

Note that we cannot really test for "instrument exogeneity" assumption and need to consider economic behavior in order to maintain the  $Cov(z, \epsilon) = 0$  assumption. At times, there may be an observable proxy for some factor contained in  $\epsilon$  and we can check if  $z$  and the proxy variable are more or less uncorrelated. And, of course, as discussed above, if we have a good proxy then we would add that variable to the equation and estimate the expanded form by OLS.

This is exactly where we see a tension between a good proxy vs a good instrumental variable. For *IQ* to be a good proxy, it needs to be as highly correlated with *ability* as possible. Yet for *IQ* to be a good instrumental variable, it needs to be uncorrelated with *ability* since *ability* is contained in  $\epsilon$  and a good instrumental variable should not covary with the error term. That is, a good instrumental variable should affect  $\log(wage)$  only through its influence on *educ* and not in any other way.

Thus, in this question, although *IQ* is a good candidate as a proxy variable for *ability*, it is not a good instrumental variable for *educ*.

---

## QUESTION 3

The following regression explores the relationship between television watching and childhood obesity, using a cross-section of US children. The variables are:

Name	Description	Minimum	Maximum	Mean
tvyst	hours of TV watched yesterday	0	6	3.14
black	dummy, 1 if black	0	1	0.31
hisp	dummy, 1 if hispanic	0	1	0.36
ageyrs	age in years	5	16	9.4
bmi	child's Body Mass Index	11	55	19
dadbmi	father's BMI	11	58	26
mombmi	mother's BMI	14	56	26

The output from a 2SLS regression appears below:

```
Instrumental-variables 2SLS regression      Number of obs   =      4,922
                                           Wald chi2(4)    =      164.47
                                           Prob > chi2     =      0.0000
                                           R-squared      =      0.0365
                                           Root MSE      =      1.7619
```

tvyst	Coefficient	Std. err.	z	P> z	[95% conf. interval]
bmi	.0452991	.0210727	2.15	0.032	.0039973 .0866009
black	.7325407	.0626985	11.68	0.000	.6096538 .8554276
hisp	.4023531	.0638145	6.31	0.000	.2772791 .5274272
ageyrs	-.0280529	.0163226	-1.72	0.086	-.0600446 .0039387
_cons	2.178131	.2608921	8.35	0.000	1.666792 2.68947

Endogenous: bmi

Exogenous: black hisp ageyrs dadbmi mombmi

Now answer the following questions.

(a) Why might an OLS regression of *tvyst* on the child's BMI give us inconsistent estimates of the causal effect of BMI on TV watching?

**Answer:** Recall that correlation between the error term and any of the regressors generally causes all of the OLS estimators to be inconsistent. In fact, if the error term is correlated with any of the independent variables, then OLS is both biased and inconsistent. This means any bias persists even as the sample size grows.

Here, if we only regress *tvyst* on *bmi* then inevitably all the omitted variables would be contained in the error term and they would be correlated with *bmi*, which would give us inconsistent estimates of the causal effect of *bmi* on tv watching.

(b) Interpret the coefficient 0.73 on *black*.

**Answer:** The coefficient implies that holding other variables constant, black children watched on average about 0.73 hours more tv than non-black children.

---

(c) Can you state a reason why we may doubt the validity of the 2SLS estimates reported above?

**Answer:** In the least, the 2SLS estimation method have the following assumptions:

- the error term of the structural equation is uncorrelated with each of the exogenous explanatory variables
- there exists at least one exogenous variable that is partially correlated with the endogenous variable in the structural equation but itself is not in the structural equation to ensure consistency
- the structural error term cannot depend on any of the exogeneous variables, i.e. homoskedasticity assumption. This ensures the 2SLS standard errors and *t*-statistics to be asymptotically valid.

Violation of any one of these assumptions would make us doubt the validity of the 2SLS estimates reported above.

## SUPPLEMENTARY QUESTIONS

### QUESTION 1

Consider the simple regression model:

$$Y_i = \alpha + \beta X_i + \varepsilon_i, \quad i = 1, 2, \dots, m \quad (1)$$

where  $Y_i$  is the mean expenditure on alcohol in group  $i$  and  $X_i$  is the mean income of group  $i$ . Each group  $i$  has  $N_i$  members and the model satisfies all the classical assumptions except that the variance of  $\varepsilon_i$  is equal to  $\sigma^2/N_i$ .

(a) What are the statistical properties of the OLS estimates of  $\alpha$  and  $\beta$  in this case?

**Answer:** Recall that when demonstrating unbiasedness and consistency of OLS estimators, homoskedasticity assumption did not play any role. That is, if the variance of the unobserved error is not constant, i.e. heteroskedastic, it does not impact whether an estimator is unbiased or consistent. Similarly, the interpretation of the goodness-of-fit measures,  $R^2$  and  $\bar{R}^2$ , are also unaffected by the presence of heteroskedasticity.

The problem with the presence of heteroskedasticity is that the estimators of the variances are biased. Since the OLS standard errors are based on these variances, they are no longer valid for constructing confidence intervals and  $t$ -statistics. In this situation the OLS  $t$ -statistics do not have  $t$  distributions and the problem is not resolved by increasing the sample size. Similarly,  $F$ -statistics are not longer  $F$ -distributed. Finally, the OLS is no longer BLUE as it is no longer asymptotically efficient.

Recall that the OLS estimator is

$$\hat{\beta} = \beta + \frac{\sum_{i=1}^m (X_i - \bar{X})\varepsilon_i}{\sum_{i=1}^m (X_i - \bar{X})^2} = \beta + \frac{\sum_{i=1}^m (X_i - \bar{X})\varepsilon_i}{SST_X^2}$$

and its variance when homoskedasticity is present is

$$\begin{aligned} Var(\hat{\beta}) &= Var\left(\beta + \frac{\sum_{i=1}^m (X_i - \bar{X})\varepsilon_i}{\sum_{i=1}^m (X_i - \bar{X})^2}\right) \\ &= Var\left(\frac{\sum_{i=1}^m (X_i - \bar{X})\varepsilon_i}{\sum_{i=1}^m (X_i - \bar{X})^2}\right) \quad \text{since } \beta \text{ is constant} \\ &= \left(\frac{1}{\sum_{i=1}^m (X_i - \bar{X})^2}\right)^2 Var\left(\sum_{i=1}^m (X_i - \bar{X})\varepsilon_i\right) \quad \text{since we are conditioning on } X_i, SST_X \text{ is nonrandom} \end{aligned}$$

$$\begin{aligned}
&= \left( \frac{1}{SST_X} \right)^2 \left( \sum_{i=1}^m (X_i - \bar{X})^2 \text{Var}(\varepsilon_i) \right) \quad \text{since we are conditioning on } X_i, X_i - \bar{X} \text{ is nonrandom} \\
&= \left( \frac{1}{SST_X} \right)^2 \left( \sum_{i=1}^m (X_i - \bar{X})^2 \sigma_\varepsilon^2 \right) \quad \text{since } \text{Var}(\varepsilon_i) = \sigma_\varepsilon^2 \text{ for all } i \text{ when homoskedastic} \\
&= \sigma_\varepsilon^2 \left( \frac{1}{SST_X} \right)^2 SST_X \\
&= \frac{\sigma_\varepsilon^2}{SST_X} = \frac{\sigma_\varepsilon^2}{\sum_{i=1}^m (X_i - \bar{X})^2}
\end{aligned}$$

and its variance when heteroskedasticity is present is

$$\text{Var}(\hat{\beta}) = \left( \frac{1}{SST_X} \right)^2 \left( \sum_{i=1}^m (X_i - \bar{X})^2 \sigma_i^2 \right) = \frac{\sum_{i=1}^m (X_i - \bar{X})^2 \sigma_i^2}{\left( \sum_{i=1}^m (X_i - \bar{X})^2 \right)^2}.$$

### Spherical Errors

We assume homoskedasticity and no autocorrelation in estimating the variance of OLS estimates. That is, we assume that all errors have the same variance  $\sigma^2$  and that there is no correlation across errors. If these hold true, then we have *spherical errors*, or that the error term follows a *spherical distribution*. This is represented in matrix form as follows:

$$\mathbb{E}(\vec{u}\vec{u}^T | \mathbf{X}) = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix} = \sigma^2 \mathbf{I}$$

To see why this is called a spherical distribution lets look at a special case of two dimensions, i.e. circular distribution, as opposed to three dimensions for spherical distribution. Consider two random errors,  $u_i$  and  $u_j$  which are graphed below as density plots and contour plots, the latter of which shows what you'd see when you look straight down from the top of the density plot.

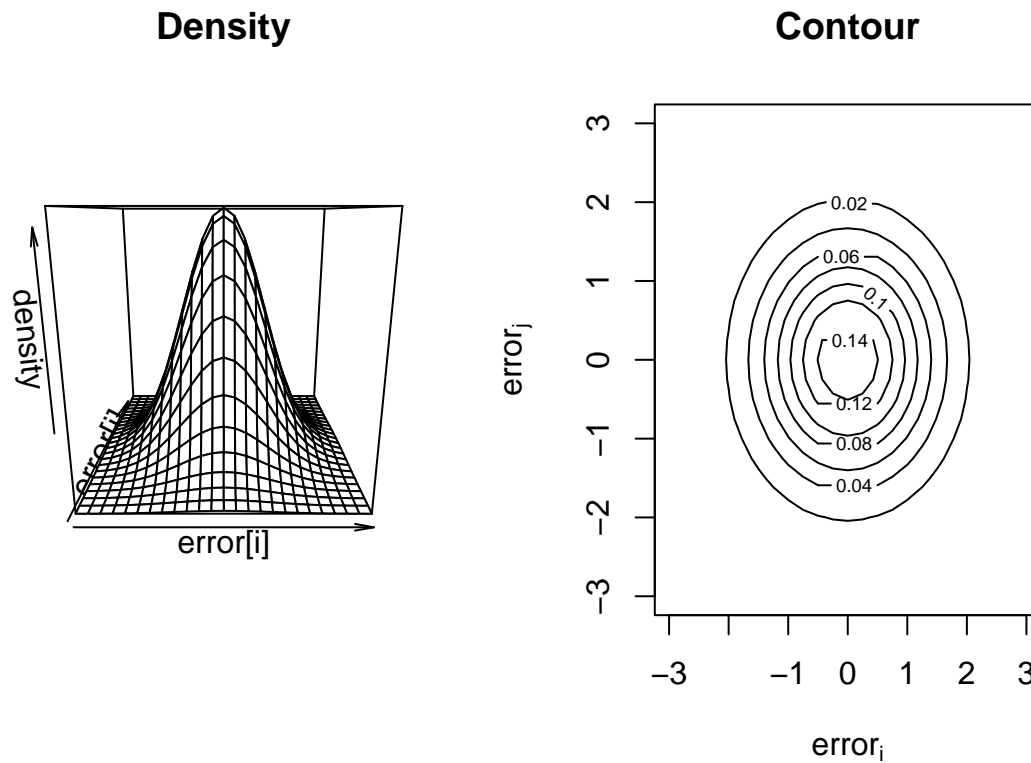
The shapes of these plots depend on the variances and covariances of these two random errors. If  $u_i$  and  $u_j$  are homoskedastic and they are not correlated, then the contour lines will be circles. If there were three random error variables  $u_i$ ,  $u_j$ , and  $u_k$  then we would have four-dimensional density plot and the contours would form a sphere. If there were more than three random error variables then the contours would form a hyper-sphere. This is why the errors are spherically distributed.

What we are plotting is therefore:

$$\mathbb{E} \begin{pmatrix} u_i \\ u_j \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad ; \quad \text{Var} \begin{pmatrix} u_i \\ u_j \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

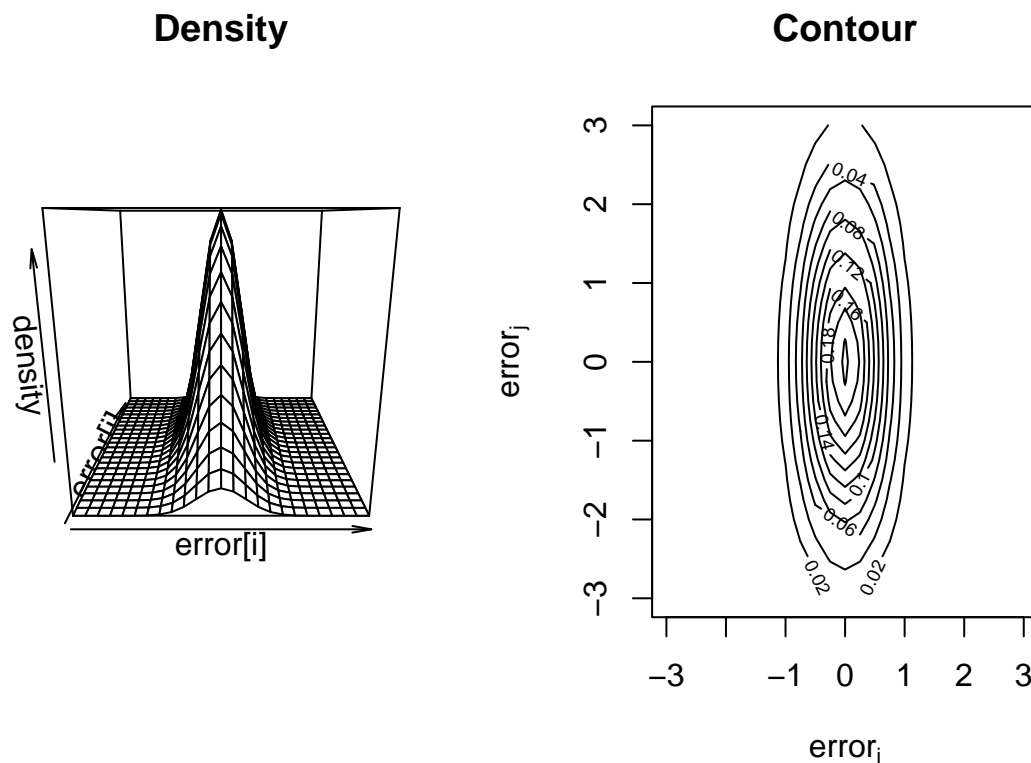
where errors are homoskedastic and there is no autocorrelation.

\end{description}



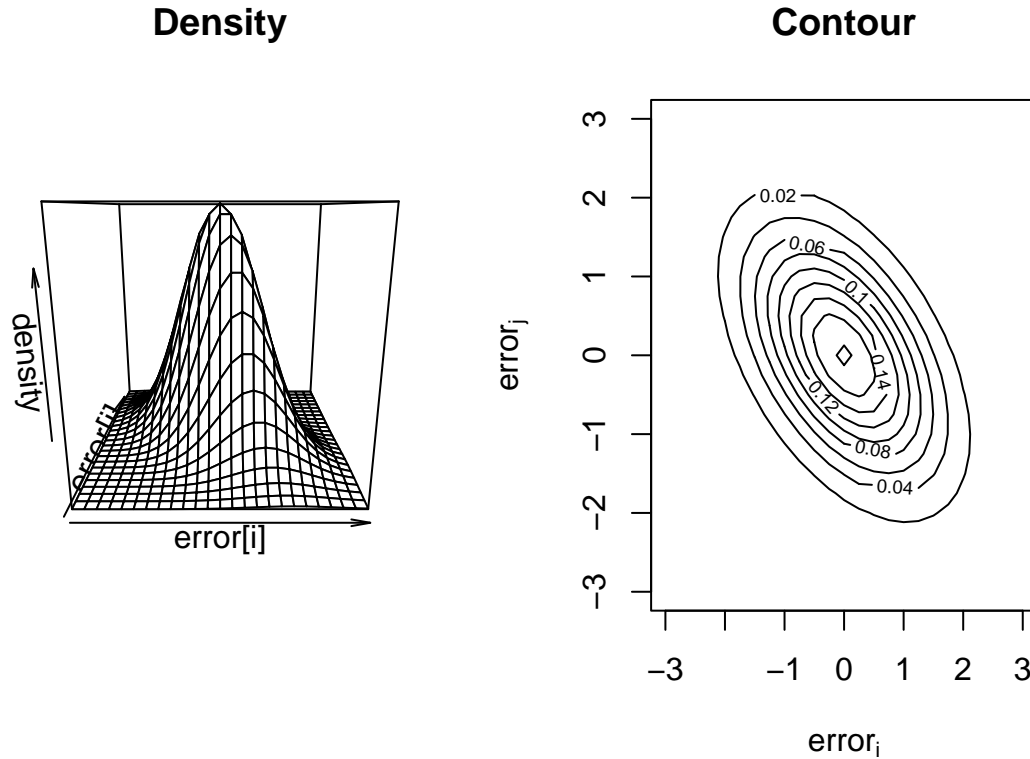
If on the other hand, heteroskedasticity is present then we lose the symmetry of the joint density plot and get a more elliptic contours. Suppose now the variance-covariance matrix is as follows:

$$\text{Var} \begin{pmatrix} u_i \\ u_j \end{pmatrix} = \begin{pmatrix} 0.25 & 0 \\ 0 & 2 \end{pmatrix}$$



Similarly, we would also get an elliptic contours if the errors are homoskedastic but there is autocorrelation. The slope of the main axis of the ellipse would depend on the sign of the correlation between the errors. Suppose now the variance-covariance matrix is as follows:

$$\text{Var} \begin{pmatrix} u_i \\ u_j \end{pmatrix} = \begin{pmatrix} 1 & -0.5 \\ -0.5 & 1 \end{pmatrix}$$



(b) How should equation (1) on page 12 be transformed so that the OLS estimates of  $\alpha$  and  $\beta$  are BLUE?

**Answer:** The variances of the error terms are given in the question, thus *known*. We can therefore estimate using the generalized least squares (GLS) estimators for correcting heteroskedasticity where we minimize a *weighted sum of squared residuals*.

↔ For remedial measures when  $\sigma_i^2$  is unknown, see Question 2 below.

$$\text{Var}(\varepsilon_i) = \frac{\sigma^2}{N_i} = \sigma_i^2$$

So we transform equation (1) on page 12 by dividing it with theses known standard deviations,  $\sigma_i$ :

$$\frac{Y_i}{\sigma_i} = \frac{\alpha}{\sigma_i} + \beta \frac{X_i}{\sigma_i} + \frac{\varepsilon_i}{\sigma_i}$$

so that

$$\begin{aligned}
 Var\left(\frac{\varepsilon_i}{\sigma_i}\right) &= \mathbb{E}\left[\left(\frac{\varepsilon_i}{\sigma_i}\right)^2\right] - \left[\mathbb{E}\left(\frac{\varepsilon_i}{\sigma_i}\right)\right]^2 \\
 &= \mathbb{E}\left[\left(\frac{\varepsilon_i}{\sigma_i}\right)^2\right] \quad \text{since } \mathbb{E}\left(\frac{\varepsilon_i}{\sigma_i}\right) = 0 \\
 &= \frac{1}{\sigma_i^2} \mathbb{E}(\varepsilon_i^2) \quad \text{since } \sigma_i^2 \text{ is known, thus it is a collection of constants} \\
 &= \frac{1}{\sigma_i^2} \sigma_i^2 = 1
 \end{aligned}$$

which is a constant. This means, the variance of the transformed disturbance term  $\frac{\varepsilon_i}{\sigma_i}$  is now homoskedastic. Since all the other assumptions of classical model still hold true, this means that if we apply OLS method to the transformed model, we will get estimators that are BLUE.

Thus, GLS is OLS on the transformed variables that satisfy the standard least-squares assumptions. The estimators that are obtained these way are GLS estimators which are BLUE.

(c) Derive  $\hat{\alpha}$  in terms of  $\hat{\beta}$  in this case.

**Answer:** In this case, what we want is a transformation of the equation 1 on page 12 in such a way that the variance of the transformed error,  $Var(\varepsilon_i^*)$ , is constant  $\sigma^2$ .

For this, we can work backwards. We know that  $Var(\varepsilon_i) = \mathbb{E}(\varepsilon_i^2) = \frac{\sigma^2}{N_i}$  so if the transformation resulted in  $Var(\varepsilon_i^*) = N_i \mathbb{E}(\varepsilon_i^2)$  then it would equal to constant  $\sigma^2$ . For that to happen, we can set  $\varepsilon_i^* = \varepsilon_i \sqrt{N_i}$ , so that

$$Var(\varepsilon_i^*) = \mathbb{E}((\varepsilon_i^*)^2) - [\mathbb{E}(\varepsilon_i^*)]^2 = \mathbb{E}((\varepsilon_i^*)^2) = \mathbb{E}((\varepsilon_i \sqrt{N_i})^2) = N_i \mathbb{E}(\varepsilon_i^2) = N_i \frac{\sigma^2}{N_i} = \sigma^2$$

as desired.

Thus using the weighting of  $\sqrt{N_i}$  the sample regression function becomes:

$$\begin{aligned}
 Y_i \sqrt{N_i} &= \alpha \sqrt{N_i} + \beta \sqrt{N_i} X_i + \varepsilon_i \sqrt{N_i} \\
 Y_i^* &= \alpha^* + \beta^* X_i + \varepsilon^*
 \end{aligned}$$

In general, to obtain the estimators for the coefficients, the weighted least-squares method minimizes the weighted residual sum of squares:

$$\sum w_i \hat{\varepsilon}_i^2 = \sum w_i (Y_i - \hat{\alpha}^* - \hat{\beta}^* X_i)^2$$

where  $\alpha^*$  and  $\beta^*$  are the weighted least squares estimators. Differentiating these with respect to  $\hat{\alpha}^*$  and  $\hat{\beta}^*$  gives us:

$$\frac{\partial}{\partial \hat{\alpha}^*} \sum w_i \hat{\varepsilon}_i^2 = 2 \sum w_i (Y_i - \hat{\alpha}^* - \hat{\beta}^* X_i) (-1)$$



$$\frac{\partial}{\partial \hat{\beta}^*} \sum w_i \hat{\varepsilon}_i^2 = 2 \sum w_i (Y_i - \hat{\alpha}^* - \hat{\beta}^* X_i)(-X_i)$$

Setting these equal to 0 gives us:

$$\begin{aligned} \sum w_i Y_i &= \hat{\alpha}^* \sum w_i + \hat{\beta}^* \sum w_i X_i \\ \sum w_i X_i Y_i &= \hat{\alpha}^* \sum w_i X_i + \hat{\beta}^* \sum w_i X_i^2 \end{aligned}$$

Solving these simultaneously, we get:

$$\begin{aligned} \hat{\alpha}^* &= \frac{\sum w_i Y_i}{\sum w_i} - \hat{\beta}^* \frac{\sum w_i X_i}{\sum w_i} \\ &= \bar{Y}^* - \hat{\beta}^* \bar{X}^* \\ \hat{\beta}^* &= \frac{(\sum w_i)(\sum w_i X_i Y_i) - (\sum w_i X_i)(\sum w_i Y_i)}{(\sum w_i)(\sum w_i X_i^2) - (\sum w_i X_i)^2} \end{aligned}$$

Notice that in this question  $w_i = N_i$  and not  $\sqrt{N_i}$ .

-----

## QUESTION 2

Using the Heteroskedasticity worksheet in sup4.xls

Load the data in R:

```
property_df <- read_excel("../Data/sup4.xls")

# You can use any of the following to examine data frame (df):
# `dim()`: for its dimensions, by row and column
# `str()`: for its structure
# `summary()`: for summary statistics on its columns
# `colnames()`: for the name of each column
# `head()`: for the first 6 rows of the data frame
# `tail()`: for the last 6 rows of the data frame
# `View()`: for a spreadsheet-like display of the entire data frame
```

(a) Estimate the following and comment on your results:

$$PRICE_t = \beta_0 + \beta_1 LOTSIZE_t + \beta_2 SQRFT_t + \beta_3 BDRMS_t + u_t \quad (2)$$

In R run the following:

```
SQ2a_lm <- lm(PRICE ~ BDRMS + LOTSIZE + SQRFT, data = property_df)
print(summary(SQ2a_lm), digits=7)
```

and in STATA run the following:

```
/* load the data */
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("heteroscedasticity") firstrow

/* `firstrow` indicates that the first row contains the variable names */
/* `describe` command would give basic information about the data set */

/* run the regression */
regress PRICE LOTSIZE SQRFT BDRMS
```

Source	SS	df	MS	Number of obs	=	88
Model	6.1713e+11	3	2.0571e+11	F(3, 84)	=	57.46
Residual	3.0072e+11	84	3.5800e+09	Prob > F	=	0.0000
				R-squared	=	0.6724
				Adj R-squared	=	0.6607
Total	9.1785e+11	87	1.0550e+10	Root MSE	=	59833

PRICE	Coefficient	Std. err.	t	P> t	[95% conf. interval]
LOTSIZE	2.067707	.6421258	3.22	0.002	.790769 3.344644
SQRFT	122.7782	13.23741	9.28	0.000	96.45415 149.1022
BDRMS	13852.52	9010.145	1.54	0.128	-4065.14 31770.18
_cons	-21770.31	29475.04	-0.74	0.462	-80384.66 36844.04

We see that the  $F$ -stat is high at 57.46 with its  $p$  value being 0. We also see that both  $LOTSIZE$  and  $SQRFT$  are significant with  $t$ -values 3.22 and 9.28 with near 0, or 0,  $p$ -values, respectively. On the other hand,  $BDRMS$  look insignificant with  $t$ -value at 1.54, though it may perhaps be due to multicollinearity.

To check for heteroskedasticity, usually the first thing to do is to plot the residuals against the estimated values of the independent variable as an amalgamation of all the dependent variables.

In R we do this with the following:

```
# the following will provide four important plots that are usually needed
# since there are four graphs, we want to display in 2x2 format first then plot
par(mfrow = c(2,2))
plot(SQ2a_lm)

# if it is only the residuals vs fitted that we are interested, then
plot(SQ2a_lm, which=1)
# or
plot(fitted(SQ2a_lm), resid(SQ2a_lm))
# we can also add a horizontal line at 0
abline(0,0)

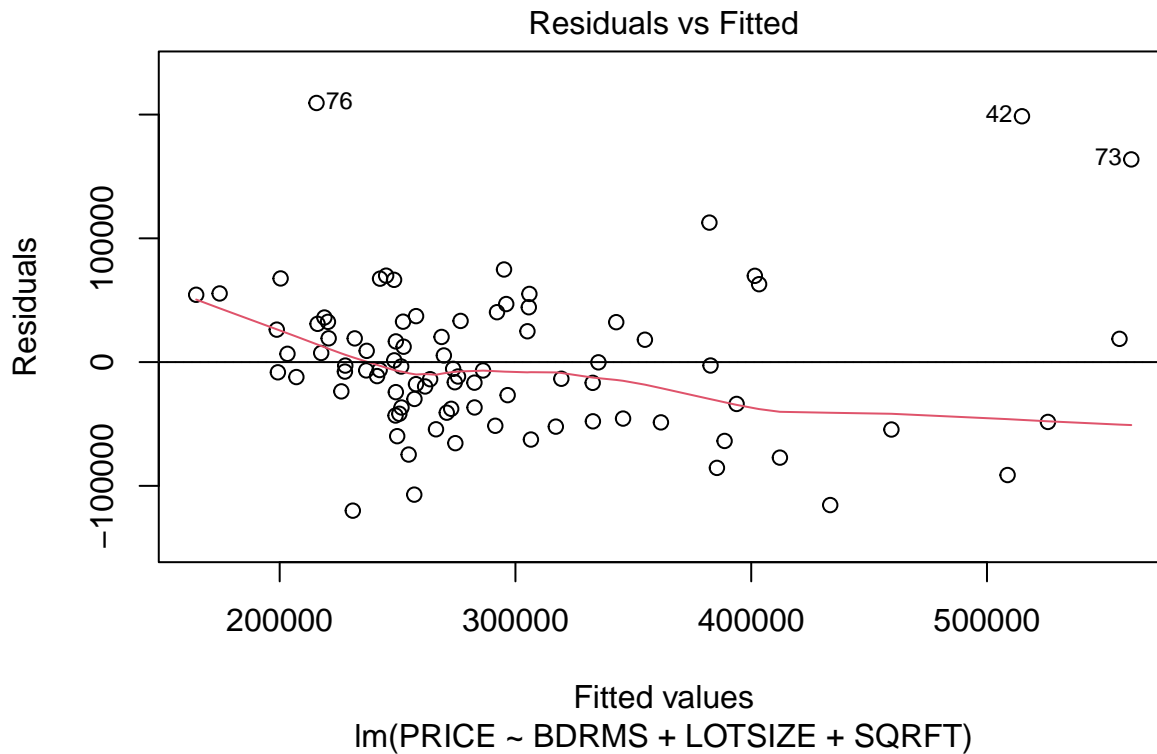
# to make this look nicer, we can also use `autoplot` command from `ggfortify` library
```

```
library(ggfortify)
autoplot(SQ2a_lm)
```

In STATA we instead use the following:

```
/* plot residuals against fitted values */
rvfplot, yline(0)
```

In either case we get the following plot:



There seems to be a downward trend which can suggest heteroskedasticity but it is difficult to tell, as it could be due to outliers.

**(b) Calculate robust standard errors for the equation 2 specified on page 17 and compare your results.**

**Answer:** White (1980)<sup>4</sup> has shown that asymptotically consistent estimates of variances and covariances of OLS estimators can be obtained even if there is heteroskedasticity present so that asymptotically valid

<sup>4</sup>White, H (1980) "A Heteroscedasticity Consistent Covariance Matrix Estimator and a Direct Test of Heteroscedasticity", *Econometrica*, 48:817-828. Though the possibility of such heteroskedasticity-robust standard errors were previously discussed by Eicker (1967) and Huber (1967) and so sometimes these are also called *White-Huber-Eicker standard errors*. See Eicker, F (1967) "Limit Theorems for Regressions with Unequal and Dependent Errors", *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability* 1:59-82, and Huber, P J (1967) "The Behavior of Maximum Likelihood Estimates under Nonstandard Conditions", *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability* 1:221-233.

statistical inferences can be made about the true parameter values. White's heteroskedasticity-corrected standard errors are also known as *robust standard errors*.

#### White's robust standard errors

##### How do we get heteroskedasticity-consistent variances and standard errors?<sup>a</sup>

Consider the simple regression model:

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

where  $Var(u_i) = \sigma_i^2$ ; that is, it is heteroskedastic. In Question 1 part (a) we have shown that

$$Var(\hat{\beta}_1) = \frac{\sum_{i=1}^n (X_i - \bar{X})^2 \sigma_i^2}{\left( \sum_{i=1}^n (X_i - \bar{X})^2 \right)^2} \quad (3)$$

Since  $\sigma_i^2$  are not directly observable, White argues for using the squared residual of each  $i$ ,  $\hat{u}_i^2$ , instead and estimating the variance of the estimator via:

$$Var(\hat{\beta}_1) = \frac{\sum_{i=1}^n (X_i - \bar{X})^2 \hat{u}_i^2}{\left( \sum_{i=1}^n (X_i - \bar{X})^2 \right)^2} \quad (4)$$

White has shown that when this equation 4 is multiplied by the sample size  $n$ , it converges in probability to  $\frac{\mathbb{E}[(X_i - \mu_x)^2 u_i^2]}{(\sigma_x^2)^2}$  which is the probability limit of equation 3 multiplied by  $n$ , and where  $\mu_x$  is the expected value of  $X$ , and  $\sigma_x^2$  is the population variance of  $X$ . Thus, the law of large numbers and the central limit theorem are key in establishing these convergences, which are necessary for justifying the use of standard errors to construct confidence intervals and  $t$ -statistics.

↔ One can first obtain the residuals from the usual OLS regression and then calculate the variance using equation 4. Statistical software do this automatically.

This can be extended to  $k$ -variable regression model

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_k X_{ki} + u_i$$

The variance of any partial regression coefficient, say  $\hat{\beta}_j$  is then obtained via

$$Var(\hat{\beta}_j) = \frac{\sum_{i=1}^n \hat{v}_{ji}^2 \hat{u}_i^2}{\left( \sum_{i=1}^n \hat{v}_{ji}^2 \right)^2} = \frac{\sum_{i=1}^n \hat{v}_{ji}^2 \hat{u}_i^2}{RSS_j^2} \quad (5)$$

where  $\hat{v}_{ji}$  denotes the  $i^{th}$  residual from regressing  $X_j$  on all other independent variables, and  $RSS_j$  is the residual sum of squares from this regression.

The square root of this expression in equation 5 is called **heteroskedasticity-robust standard error** for  $\hat{\beta}_j$ .

Also note that, sometimes the equation 5 is adjusted for degrees of freedom by multiplying it with  $\frac{n}{n-(k+1)}$  before taking the square root. This is because if  $\hat{u}_i$  were the same for all  $i$  then we would get the usual OLS standard errors. Since all forms of this equation has asymptotic justification, and they are asymptotically equivalent, no one form is unanimously preferred over others. Usually, we use whatever form the software we work with uses.

<sup>a</sup>Gujarati and Porter (2009), Appendix 11A.4; Wooldridge (2021), Section 8.2

We can now calculate this in R as follows:

```
#we need two additional libraries for this:
library(lmtest) #for `coeftest` function
library(sandwich) #for `vcovHC` function

coeftest(SQ2a_lm, vcov = vcovHC(SQ2a_lm, "HC1"))
#the default in vcovHC is "HC3" but to get the exact result as STATA we use "HC1"
```

Similarly, we can calculate this in STATA as follows:

```
/* run the regression with additional `robust` command */

regress PRICE LOTSIZE SQRFT BDRMS, robust
```

However, to present the “robust” and “nonrobust” results side by side in a table, we can use the following set of commands instead:

```
/* load the data */
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("heteroscedasticity") firstrow

/* run the regression with `robust` command */
quietly regress PRICE LOTSIZE SQRFT BDRMS, robust

/* store the estimates under the heading "robust" */
estimates store robust

/* run the regression for nonrobust */
quietly regress PRICE LOTSIZE SQRFT BDRMS

/* store the estimates under the heading "nonrobust" */
estimates store nonrobust

/* create the table for robust and nonrobust estimates of beta, s.e., and t-values */
estimates table robust nonrobust, b se t
```

Variable	robust	nonrobust
LOTSIZE	2.0677066	2.0677066
	1.2514244	.64212582
	1.65	3.22
SQRFT	122.77819	122.77819
	17.725334	13.237407
	6.93	9.28
BDRMS	13852.522	13852.522
	8478.625	9010.1454
	1.63	1.54
_cons	-21770.309	-21770.309
	37138.211	29475.042
	-0.59	-0.74

Legend: b/se/t

Notice that all the  $t$ -values are lower for each variable and in the case of *LOTSIZE* this reduction means it is no longer significant. This is at least suggestive that *LOTSIZE* may be a major source of the heteroskedasticity.

(c) Using the specification in equation (2) on page 17, conduct a Goldfeld-Quandt test for heteroskedasticity in the *LOTSIZE* dimension (exclude the middle 24 observations).

**Answer:** The Goldfeld-Quandt<sup>5</sup> test is applicable when we assume that the heteroskedastic variance,  $\sigma_i^2$ , is positively related to *one* of the explanatory variables in the regression model. In this question we are assuming that the heteroskedastic variance is related to *LOTSIZE*.

#### Goldfeld-Quandt Test

**What are the reasoning and mechanics of the test?<sup>a</sup>**

Consider the simple regression model:

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

and suppose  $\sigma_i^2$  is monotonically related to  $X_i$ . One plausible assumption of this is

$$\sigma_i^2 = \sigma^2 X_i^2. \quad (6)$$

What this assumption says is that  $\sigma_i^2$  is proportional to the square of the  $X$  variable. If this assumption is appropriate, it would mean the larger  $X_i$  values are, the larger  $\sigma_i^2$  gets. If that turns out to be the case, heteroskedasticity is most likely to be present in the model.

To test this, Goldfeld and Quandt provide the following steps:

Step 1: Order or rank the observations according to the values of  $X_i$  beginning with the lowest  $X$  value;

Step 2: Omit  $c$  central observations, where  $c$  is specified a priori, and divide the remaining observations into two groups, each of  $\frac{n-c}{2}$  observations;

Step 3: Fit separate OLS regressions to these two groups of observations and obtain the respective residual sum of squares  $RSS_1$  and  $RSS_2$ , where  $RSS_1$  represents the  $RSS$  from the regression corresponding to the smaller  $X_i$  values, i.e. small variance group, and  $RSS_2$  to the larger  $X_i$  values, i.e. the large variance group.

These  $RSS$  each have  $\frac{n-c}{2} - (k+1)$  degrees of freedom where  $k$  is the number of parameters to be estimated, excluding the intercept - hence  $+1$ .

Step 4: Compute the following ratio:

$$\lambda = \frac{\frac{RSS_2}{df}}{\frac{RSS_1}{df}}$$

The main argument of this test is that if the assumption of homoskedasticity and  $u_i$  are normally distributed both hold true, then  $\lambda$  of equation (6) follows the  $F$ -distribution with  $\frac{n-c}{2} - (k+1)$  degrees of freedom in both the numerator and denominator.

<sup>5</sup>Goldfeld S, and Quandt R E (1972) *Nonlinear Methods in Econometrics*, North Holland Publishing Company, Amsterdam.

As usual, if the computed  $\lambda$  which is equal to  $F$ -statistic, is greater than the critical  $F$  value at the chosen level of significance, we can reject the null hypothesis of homoskedasticity.

**Why we omit  $c$  central observations?** These observations are omitted to accentuate the difference between the small variance group,  $RSS_1$ , and the large variance group  $RSS_2$ . However, the *power* of the test depends on how  $c$  is chosen. Recall that *power of a test* is measured by the probability of rejecting the null hypothesis when it is false, and it is calculated by  $1 - \text{prob}(\text{Type II error})$ .

Goldfeld and Quandt suggest  $c = 8$  for models with two-explanatory variables if  $n = 30$  and double if  $n = 60$ .

<sup>a</sup>Gujarati and Porter (2009), Section 11.5

In this question  $c = 24$  and we order *LOTSIZE* from small to large. To run the Goldfeld-Quandt test in R we can use the `gqtest()` function from the `lmtest` library:

```
gqtest(SQ2a_lm, order.by = property_df$LOTSIZE, fraction = 24, alternative="two.sided")
```

Goldfeld-Quandt test

```
data: SQ2a_lm
GQ = 1.6275, df1 = 28, df2 = 28, p-value = 0.2037
alternative hypothesis: variance changes from segment 1 to 2
```

```
qf(0.975, 28, 28, lower.tail = TRUE) #critical F-value
```

```
[1] 2.129924
```

In STATA there are more steps involved. First we need to order the data and removed the middle 24 observations before running regression on each:

```
/* load the data */
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("heteroscedasticity") firstrow

/* Step 1: Order the data according to LOTSIZE values */
sort LOTSIZE

/* create an index on which we will impose our condition for splitting data */
gen index=_n

/* run the regressions on each splitted data */
reg PRICE BDRMS LOTSIZE SQRFT if index<33

reg PRICE BDRMS LOTSIZE SQRFT if index>56

/* Derive F-stat by dividing RSS of each (since df cancel out) */
display e(rss)/8.5839e+10

/* compute critical F value */
display invfprob(28,28,0.025)
```

Source	SS	df	MS	Number of obs	=	32
				F(3, 28)	=	0.95
Model	8.7445e+09	3	2.9148e+09	Prob > F	=	0.4295
Residual	8.5839e+10	28	3.0657e+09	R-squared	=	0.0925
				Adj R-squared	=	-0.0048
Total	9.4584e+10	31	3.0511e+09	Root MSE	=	55369

PRICE	Coefficient	Std. err.	t	P> t	[95% conf. interval]	
BDRMS	2118.224	14818.04	0.14	0.887	-28235.16	32471.61
LOTSIZE	6.442439	9.678	0.67	0.511	-13.38204	26.26692
SQRFT	55.13037	32.93973	1.67	0.105	-12.34362	122.6044
_cons	104434.3	100019.7	1.04	0.305	-100446.7	309315.4

Source	SS	df	MS	Number of obs	=	32
				F(3, 28)	=	24.00
Model	3.5918e+11	3	1.1973e+11	Prob > F	=	0.0000
Residual	1.3970e+11	28	4.9894e+09	R-squared	=	0.7200
				Adj R-squared	=	0.6900
Total	4.9888e+11	31	1.6093e+10	Root MSE	=	70636

PRICE	Coefficient	Std. err.	t	P> t	[95% conf. interval]	
BDRMS	22392.55	16559.13	1.35	0.187	-11527.28	56312.39
LOTSIZE	1.385209	.8391777	1.65	0.110	-.3337691	3.104186
SQRFT	121.5091	23.75068	5.12	0.000	72.85808	170.1602
_cons	-26511.29	50274.72	-0.53	0.602	-129494.4	76471.81

1.6275086

2.1299243

Both R and STATA give the same result that the  $F$ -statistic of 1.6275654 is smaller than the critical  $F$ -value of 2.1299243 which means we cannot reject the null of homoskedasticity. Therefore, it seems as though there is no heteroskedasticity according to the Goldfeld-Quandt test. However, the form of heteroskedasticity may be more complicated.

(d) Test for heteroskedasticity by first estimating an equation that regresses the squared residuals from equation (2) on page 17 against all of the independent variables used to estimate equation (2). (Calculate both F and LM versions of this test). Verify your results using the 'hettest' command in Stata. Compare these results with the results of the White Test in Stata.

**Answer:** Goldfeld-Quandt test depends not only on the number of observations we omit but also on identifying the correct  $X$  variable that needs to be ordered. These limitations of this test can be avoided



with *Breusch-Pagan Test*,<sup>6</sup> or BP test, which is also called *Breusch-Pagan-Godfrey Test*,<sup>7</sup> or BPG test.

#### Breusch-Pagan / Breusch-Pagan-Godfrey Test

##### What are the reasoning and mechanics of the test?<sup>a</sup>

Consider the model

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_k X_k + u$$

and assume that  $\mathbb{E}(u|X_1, X_2, \dots, X_k) = 0$  so that OLS is unbiased and consistent.

The null hypothesis is that homoskedasticity holds and we require the data to tell us otherwise. That is,  $\mathbb{H}_0 : \text{Var}(u|X_1, X_2, \dots, X_k) = \sigma^2$ ; and since  $\text{Var}(u|X_1, X_2, \dots, X_k) = \mathbb{E}(u^2|X_1, X_2, \dots, X_k)$  the null hypothesis can be expressed as:

$$\mathbb{H}_0 : \text{Var}(u|X_1, X_2, \dots, X_k) = \mathbb{E}(u^2|X_1, X_2, \dots, X_k) = \mathbb{E}(u^2) = \sigma^2.$$

This shows that in order to test for violation of the homoskedasticity assumption we want to test whether  $u^2$  is related in expected value to one or more of the explanatory variables. Therefore, if  $\mathbb{H}_0$  is false, then the expected value of  $u^2$  given the independent variables, i.e.  $\mathbb{E}(u^2|X_1, X_2, \dots, X_k)$  can be any function of the  $X_j$ . A simple approach is to assume a linear function:

$$u^2 = \gamma_0 + \gamma_1 X_1 + \gamma_2 X_2 + \cdots + \gamma_k X_k + v$$

The null hypothesis then becomes

$$\mathbb{H}_0 : \gamma_1 = \gamma_2 = \cdots = \gamma_k = 0.$$

Under the null hypothesis we can assume that the error  $v$  is independent of  $X_1, \dots, X_k$ . Then, either the  $F$  or *Lagrange Multiplier (LM)* statistics can be used to test for the overall significance of the independent variables in explaining  $u^2$ . Both statistics would have asymptotic justification, even though  $u^2$  cannot be normally distributed.

↪ e.g. if  $u$  is normally distributed then  $\frac{u^2}{\sigma^2}$  is distributed  $\chi_1^2$ .

If we could observe the  $u^2$  in the sample, then we could compute this statistic by running the OLS regression of  $u^2$  on  $X_1, \dots, X_k$  using all  $n$  observations, which would give us the maximum likelihood (ML) of  $\sigma^2$ .

Since we do not know  $u$ , we can instead estimate the equation:

$$\hat{u}^2 = \gamma_0 + \gamma_1 X_1 + \gamma_2 X_2 + \cdots + \gamma_k X_k + \epsilon$$

and compute  $F$  or  $LM$  statistics for the joint significance of  $X_1, \dots, X_k$ . The  $F$  and  $LM$  statistic both depend on the  $R$ -squared value of this regression,  $R_{\hat{u}^2}^2$ .

The  $F$ -statistic for heteroskedasticity is

$$F = \frac{\frac{R_{\hat{u}^2}^2}{k}}{\frac{1 - R_{\hat{u}^2}^2}{n - (k + 1)}}$$

where  $k$  is the number of regressors. This  $F$  statistic has approximately an  $F_{k, n-(k+1)}$  distribution under the null hypothesis of homoskedasticity.

The  $LM$  statistic for heteroskedasticity is

$$LM = n \times R_{\hat{u}^2}^2$$

which is the  $R$ -squared of the error regression multiplied by the sample size. Under the null hypothesis,  $LM$  is distributed asymptotically as  $\chi_k^2$ .

<sup>6</sup>Breusch, T and Pagan A (1979) "A Simple Test for Heteroscedasticity and Random Coefficient Variation", *Econometrica*, 47:1287-1294.

<sup>7</sup>Godfrey L (1978) "Testing for Multiplicative Heteroscedasticity" *Journal of Econometrics*, 8:227-236.

The *LM* version of the test is called the **Breusch-Pagan test** for heteroskedasticity, or BP-test; though the *LM*-statistic form was suggested by Koenker (1981).<sup>b</sup>

### What is Lagrange Multiplier Statistic?<sup>c</sup>

Consider again the multiple regression model with  $k$  independent variables:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_k X_k + u$$

We want to test whether, say, the last  $q$  of these variables all have 0 population parameters. The null hypothesis is therefore

$$\mathbb{H}_0 : \beta_{k-q+1} = \beta_{k-q+2} = \cdots = \beta_k = 0,$$

which puts  $q$  exclusion restrictions on the model. The alternative hypothesis is that at least one of the parameters is different from 0.

The LM statistic requires the estimation of the restricted model only. So we run the regression

$$Y = \beta_0^{res} + \beta_1^{res} X_1 + \cdots + \beta_{k-q}^{res} X_{k-q} + u^{res}$$

where  $u^{res}$  indicate that the residuals are from the restricted model. Note that this is a shorthand to indicate that we obtain restricted residual for each observation in the sample, but didn't use the  $i$  subscript to avoid crowding of subscripts.

The idea is that if the omitted variables  $X_{k-q+1}$  through  $X_k$  truly have 0 population coefficients, then  $u^{res}$  should at least be approximately uncorrelated with each of these variables in the sample. In fact, it should be uncorrelated with all regressors because the omitted regressors in the restricted model are correlated with the regressors that appear in the restricted model.

This means, we run the regression of  $u^{res}$  on  $X_1, \dots, X_k$ .

↪ this is an example of *auxiliary regression* which is a regression used to compute a test statistic but whose coefficients are not of direct interest.

If the null hypothesis is true, the  $R$ -squared from this regression should be "close" to zero, subject to sampling error. This is because  $u^{res}$  will be approximately uncorrelated with all the independent variables.

What is interesting with this test is that, under the null hypothesis, the sample size multiplied by the  $R$ -squared from the auxiliary regression is distributed asymptotically as a chi-square random variable with  $q$  degrees of freedom. That is,  $n \times R_{u^2}^2 \stackrel{a}{\sim} \chi_q^2$ .

Because of its form, the *LM* statistic is also referred to as the **n-R-squared statistic**.

<sup>a</sup>Gujarati and Porter (2009), Section 11.5; Wooldridge (2021), Section 8.3

<sup>b</sup>Koenker, R (1981) "A Note on Studentizing a Test for Heteroskedasticity", *Journal of Econometrics* 17:107-112.

<sup>c</sup>Wooldridge (2021), Section 5.2a

We can obtain the BP-statistic that has a  $\chi_3^2$  distribution in R as follows:

```
bptest(SQ2a_lm)
```

In STATA, we can do the same using the **hettest** command:

```
/* load the data */
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("heteroscedasticity") firstrow

/* run the regression */
```

```
quietly regress PRICE LOTSIZE SQRFT BDRMS

/* apply hettest where rhs mean right-hand-side */
hettest, rhs fstat
hettest, rhs iid

/* manual calculation to compare the results */
predict u, r
generate U2 = u^2
quietly regress U2 LOTSIZE SQRFT BDRMS

/* display the F-statistic and LM-statistic */
display e(F)
display e(r2)*e(N)
```

Breusch-Pagan/Cook-Weisberg test for heteroskedasticity

Assumption: i.i.d. error terms

Variables: All independent variables

H0: Constant variance

$F(3, 84) = 5.34$

Prob > F = 0.0020

Breusch-Pagan/Cook-Weisberg test for heteroskedasticity

Assumption: i.i.d. error terms

Variables: All independent variables

H0: Constant variance

$\chi^2(3) = 14.09$

Prob >  $\chi^2$  = 0.0028

5.3389193

14.092385

it is often the case that  $\chi^2$ -tests have better properties but are harder to explain. So here we use the  $F$  initially to give the intuition then point out which  $\chi^2$ -tests do roughly the same things. Based on the  $p$ -values we can reject the null hypothesis. The Breusch-Pagan test suggests the presence of heteroskedasticity.

However, the BP test assumes that the form of the heteroskedasticity is linear. To try out different forms of the relations, we can use the White test.

#### White Test

##### What are the reasoning and mechanics of the test?<sup>a</sup>

Unlike Goldfeld-Quandt test, which requires reordering the observations with respect to the  $X$  variable that supposedly caused heteroskedasticity, or the BP test, which is sensitive to the normality and linearity assumptions, the general heteroskedasticity test proposed by White (1980)<sup>b</sup> does not rely

on the normality assumption.

White test uses the insight that the homoskedasticity assumption can be replaced with the weaker assumption that the squared error  $\hat{u}_i^2$  is *uncorrelated* with all the independent variables,  $X_j$ , the squares of the independent variables,  $X_j^2$ , and all the cross products,  $X_j X_h$  for  $j \neq h$ .

The test is explicitly intended to test for forms of heteroskedasticity that invalidate the usual OLS standard errors and test statistics. Consider a model with  $k = 3$  independent variables. The White test process is as follows:

Step 1: Estimate  $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + u_i$  and obtain the residuals,  $\hat{u}_i$ ;

Step 2: Obtain the  $R_{\hat{u}^2}^2$  from following auxiliary regression:

$$\hat{u}_i^2 = \gamma_0 + \gamma_1 X_{1i} + \gamma_2 X_{2i} + \gamma_3 X_{3i} + \gamma_4 X_{1i}^2 + \gamma_5 X_{2i}^2 + \gamma_6 X_{3i}^2 + \gamma_7 X_{1i} X_{2i} + \gamma_8 X_{1i} X_{3i} + \gamma_9 X_{2i} X_{3i} + \epsilon_i$$

That is, we are regressing the squared residuals from the original regression on the original variables, their squared values, and the cross products of the regressors. We can also introduce higher powers of regressors if necessary.

Step 3: Under the null hypothesis that there is no heteroskedasticity,  $n \times R_{\hat{u}^2}^2 \stackrel{a}{\sim} \chi_{df}^2$ . In this example we have 9 regressors, so  $df = 9$ .

Step 4: If the  $\chi_{df}^2$  value obtained is higher than the critical  $\chi_{df}^2$  at the chosen level of significance, then this test would suggest a presence of heteroskedasticity. If it does not exceed the critical value, then we cannot reject  $\mathbb{H}_0 : \gamma_1 = \dots = \gamma_9 = 0$ .

One final point is that this approach of White test uses many degrees of freedom. We can have a slightly different approach to White test that can conserve on degrees of freedom. To create the test, notice that the difference between White and BP tests is that the White test includes the squares and cross-products of the independent variables, whereas BP doesn't. We can preserve the spirit of the White test while conserving on degrees of freedom by using the OLS fitted values in a test for heteroskedasticity.

Recall the fitted values are defined for each observation  $i$  by

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \dots + \hat{\beta}_k X_{ki}$$

These are just linear functions of the independent variables. If we square the fitted values, we get a particular function of all the squares and cross products of the independent variables. This suggests testing for heteroskedasticity by estimating the equation

$$\hat{u}_i^2 = \eta_0 + \eta_1 \hat{Y}_i + \eta_2 \hat{Y}_i^2 + \epsilon$$

where  $\hat{Y}_i$  stand for fitted values. We can use  $F$  or  $LM$  statistic for the null hypothesis  $\mathbb{H}_0 : \eta_1 = 0, \eta_2 = 0$ . This results in two restrictions in testing the null of homoskedasticity, regardless of the number of independent variables in the original model. This can be thought of as a special case of White test.

<sup>a</sup>Gujarati and Porter (2009), Section 11.5; Wooldridge (2021), Section 8.3a

<sup>b</sup>White H (1980) "A Heteroskedasticity Consistent Covariance Matrix Estimator and a Direct Test of Heteroskedasticity", *Econometrica*, 48:817-818

We can run the White test in R manually as follows:

```
Ru2_SQ2 <- summary(lm(resid(SQ2a_lm) ~ fitted(SQ2a_lm) + I(fitted(SQ2a_lm)^2)))$r.squared
LM_SQ2 <- nrow(property_df)*Ru2_SQ2
p_value_SQ2 <- 1-pchisq(LM_SQ2,2)
p_value_SQ2
```

or we can use the `bptest()` function from `lmtest` package:

```
bptest(SQ2a_lm, ~ BDRMS + LOTSIZE + SQRFT
      + I(BDRMS)^2 + I(LOTSIZE)^2 + I(SQRFT)^2
      + BDRMS*LOTSIZE + BDRMS*SQRFT + SQRFT*LOTSIZE,
      data = property_df)

#Special case of White test that conserves on degrees of freedom
bptest(SQ2a_lm, ~ fitted(SQ2a_lm) + poly(fitted(SQ2a_lm),2))
```

```
/* load the data */
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("heteroscedasticity") firstrow

/* run the regression */
quietly regress PRICE LOTSIZE SQRFT BDRMS

/* manual calculation for White Test */
predict u, r
generate U2 = u^2
generate B2 = BDRMS^2
generate L2 = LOTSIZE^2
generate S2 = SQRFT^2
generate BL = BDRMS*LOTSIZE
generate BS = BDRMS*SQRFT
generate LS = LOTSIZE*SQRFT

quietly regress U2 BDRMS LOTSIZE SQRFT B2 L2 S2 BL BS LS

/* calculate the chi-square statistic */
display e(N)*e(r2)
```

33.731658

or we can use the `imtest`, `white` command in STATA after the original regression:

```
/* load the data */
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("heteroscedasticity") firstrow

/* run the regression */
quietly regress PRICE LOTSIZE SQRFT BDRMS

/* run the White test */
imtest, white
```

White's test

H0: Homoskedasticity

Ha: Unrestricted heteroskedasticity

```
chi2(9) = 33.73
Prob > chi2 = 0.0001
```

Cameron & Trivedi's decomposition of IM-test

Source	chi2	df	p
Heteroskedasticity	33.73	9	0.0001
Skewness	8.14	3	0.0432
Kurtosis	-163111.28	1	1.0000
Total	-163069.41	13	1.0000

In all of these approaches we obtain a chi-squared value of 33.73 with 0.001 p-value. Thus we can reject the null hypothesis of the homoskedasticity.

However, there is sure to be lots of multicollinearity here so it is difficult to tell if there is a non-linear relationship with any of the variables. We can run individual regressions and see what we can find. For example, with *LOTSIZE*:

```
/* load the data */
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("heteroscedasticity") firstrow

/* run the regression */
quietly regress PRICE LOTSIZE SQRFT BDRMS

/* obtain residual squareds and create LOTSIZE squared */
predict u, r
generate U2 = u^2
generate LOTSIZE2 = LOTSIZE^2

/* regress residuals on lotsize for nonlinearity */
regress U2 LOTSIZE LOTSIZE2
```

Source	SS	df	MS	Number of obs	=	88
Model	7.5587e+20	2	3.7793e+20	F(2, 85)	=	8.87
Residual	3.6229e+21	85	4.2622e+19	Prob > F	=	0.0003
Total	4.3787e+21	87	5.0330e+19	R-squared	=	0.1726
				Adj R-squared	=	0.1532
				Root MSE	=	6.5e+09

U2	Coefficient	Std. err.	t	P> t	[95% conf. interval]
LOTSIZE	733024.8	209062.3	3.51	0.001	317352.9 1148697
LOTSIZE2	-5.897015	2.318524	-2.54	0.013	-10.50686 -1.287167
_cons	-2.11e+09	1.64e+09	-1.28	0.203	-5.38e+09 1.16e+09

There appears to be a non-linear relationship to *LOTSIZE*, in which case the White test may be better than the Breusch-Pagan test.

(e) Researcher A decides to try ‘scaling’ to remove the heteroskedasticity and so reestimates the equation in part (a) by dividing equation (2) on page 17 by  $BDRMS$  and then, as an alternative, by  $LOTSIZE$ . Discuss the reasoning behind using such variable to scale in this way (what must the form of heteroskedasticity be in each case? - note the  $1/X$  term in each regression you have to estimate!). Which gives the best results?

**Answer:** Recall in Question 2(b) we calculated White’s robust standard errors. That method has some drawbacks, however. In addition to being a large-sample procedure, the estimators obtained using White’s robust standard errors may not be so efficient compared to those obtained by transforming the data to reflect specific types of heteroskedasticity.

To see this, consider the simple regression

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

and assume that the error variance is proportional to square of the explanatory variable,  $X_i^2$ :

$$\mathbb{E}(u_i^2) = \sigma^2 X_i^2.$$

If this is the case, the original model can be transformed to yield a homoskedastic error variance as follows.

First, divide the original model by  $X_i$  so that

$$\begin{aligned} \frac{Y_i}{X_i} &= \frac{\beta_0}{X_i} + \beta_1 \frac{X_i}{X_i} + \frac{u_i}{X_i} \\ &= \frac{\beta_0}{X_i} + \beta_1 + v_i \end{aligned}$$

where  $v_i = u_i/X_i$  is the transformed disturbance term. Then,

$$\mathbb{E}(v_i^2) = \mathbb{E}\left[\left(\frac{u_i}{X_i}\right)^2\right] = \frac{1}{X_i^2} \mathbb{E}(u_i^2) = \frac{1}{X_i^2} \sigma^2 X_i^2 = \sigma^2.$$

That is, the variance of  $v_i$  is homoskedastic and OLS can be applied to the transformed equation.

This is exactly what the question is asking us to do. In the first part we will estimate the following:

$$\frac{PRICE_i}{BDRMS_i} = \beta_0 \frac{1}{BDRMS_i} + \beta_1 \frac{LOTSIZE_i}{BDRMS_i} + \beta_2 \frac{SQRFT_i}{BDRMS_i} + \beta_3 + u_i \frac{1}{BDRMS_i}$$

In R we can run the transformed regression and run the BP and White tests as follows:

```
SQ2e_lm_bdrms <- lm(I(PRICE/BDRMS) ~ I(LOTSIZE/BDRMS) + I(SQRFT/BDRMS) + I(1/BDRMS),
  data=property_df)
summary(SQ2e_lm_bdrms)

# Breusch-Pagan test
bptest(SQ2e_lm_bdrms)

# White test
bptest(SQ2e_lm_bdrms, ~ I(LOTSIZE/BDRMS) + I(SQRFT/BDRMS) + I(1/BDRMS)
  + I((LOTSIZE/BDRMS)^2) + I((SQRFT/BDRMS)^2) + I((1/BDRMS)^2)
  + I(LOTSIZE/BDRMS)*I(SQRFT/BDRMS) + I(LOTSIZE/BDRMS)*I(1/BDRMS)
  + I(SQRFT/BDRMS)*I(1/BDRMS), data=property_df)

#Special case of White test that conserves on degrees of freedom
bptest(SQ2e_lm_bdrms, ~ fitted(SQ2e_lm_bdrms) + poly(fitted(SQ2e_lm_bdrms),2))
```

and in STATA as follows:

```
/* load the data */
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("heteroscedasticity") firstrow

/* generate the transformations */
generate PRBD = PRICE/BDRMS
generate LTBD = LOTSIZE/BDRMS
generate FTBD = SQRFT/BDRMS
generate BD = 1/BDRMS

/* run the regression */
regress PRBD LTBD FTBD BD

/* run the Breusch-Pagan test */
hettest LTBD FTBD BD, iid

/* run the White test */
imtest, white
```

Source		SS	df	MS	Number of obs	=	88
-----+-----					F(3, 84)	=	37.39
Model		2.8648e+10	3	9.5493e+09	Prob > F	=	0.0000
Residual		2.1454e+10	84	255402500	R-squared	=	0.5718
-----+-----					Adj R-squared	=	0.5565
Total		5.0102e+10	87	575882609	Root MSE	=	15981
-----+-----							
PRBD		Coefficient	Std. err.	t	P> t	[95% conf. interval]	
-----+-----							
LTBD		1.87511	.6600398	2.84	0.006	.5625488	3.187672
FTBD		107.9831	12.96118	8.33	0.000	82.20834	133.7578
BD		26623.52	27637.01	0.96	0.338	-28335.7	81582.75
_cons		8850.31	8744.765	1.01	0.314	-8539.614	26240.23
-----+-----							

Breusch-Pagan/Cook-Weisberg test for heteroskedasticity

Assumption: i.i.d. error terms

Variables: LTBD FTBD BD

H0: Constant variance

chi2(3) = 2.59  
Prob > chi2 = 0.4589

White's test

H0: Homoskedasticity

Ha: Unrestricted heteroskedasticity

chi2(9) = 14.13  
Prob > chi2 = 0.1178



Cameron & Trivedi's decomposition of IM-test

Source	chi2	df	p
Heteroskedasticity	14.13	9	0.1178
Skewness	2.57	3	0.4619
Kurtosis	-3292.58	1	1.0000
Total	-3275.87	13	1.0000

From the chi-squared values and their associated p-values we cannot reject the null hypothesis of homoskedasticity.

Note that we can also obtain the same using the GLS, or weighted least squares, approach discussed in question 1(b). Here the weights are  $1/BDRMS_i$ .

If we do the same transformation with *LOTSIZE* we get:

```
* load the data */
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("heteroscedasticity") firstrow

/* generate the transformations */
generate PRLT = PRICE/LOTSIZE
generate LT = 1/LOTSIZE
generate FTLT = SQRFT/LOTSIZE
generate BDLT = BDRMS/LOTSIZE

/* run the regression */
regress PRLT LT FTLT BDLT

/* run the Breusch-Pagan test */
hettest LT FTLT BDLT, iid

/* run the White test */
imtest, white
```

Source	SS	df	MS	Number of obs	=	88
Model	83036.7811	3	27678.927	F(3, 84)	=	480.77
Residual	4836.01679	84	57.5716285	Prob > F	=	0.0000
Total	87872.7979	87	1010.03216	R-squared	=	0.9450
				Adj R-squared	=	0.9430
				Root MSE	=	7.5876

PRLT	Coefficient	Std. err.	t	P> t	[95% conf. interval]
LT	21904.58	30415.95	0.72	0.473	-38580.88 82390.04
FTLT	97.29152	8.917066	10.91	0.000	79.55896 115.0241
BDLT	3837.126	7042.929	0.54	0.587	-10168.51 17842.76
_cons	7.358056	1.721546	4.27	0.000	3.934574 10.78154

-----

Breusch-Pagan/Cook-Weisberg test for heteroskedasticity

Assumption: i.i.d. error terms

Variables: LT FTLT BDLT

H0: Constant variance

chi2(3) = 5.96

Prob > chi2 = 0.1134

White's test

H0: Homoskedasticity

Ha: Unrestricted heteroskedasticity

chi2(9) = 8.74

Prob > chi2 = 0.4617

Cameron & Trivedi's decomposition of IM-test

Source	chi2	df	p
Heteroskedasticity	8.74	9	0.4617
Skewness	1.66	3	0.6448
Kurtosis	1.55	1	0.2126
Total	11.96	13	0.5311

Once again, from the chi-squared values and their associated p-values we cannot reject the null hypothesis of homoskedasticity.

(f) Researcher B decides to pursue the following strategies to remove heteroskedasticity:

- i) Use logged data for PRICE, LOTSIZE, SQRFRT (N.B. don't drop BDRMS)
- ii) Remove outliers (observations 42,73,76,77)

Discuss the reasoning behind each of these strategies, and the results obtained in each case.

**Answer (i):** A log transformation such as  $\ln Y_i = \beta_0 + \beta_1 \ln X_i + u_i$  often reduces heteroskedasticity because log transformation compresses the scales, reducing say a ten-fold difference between two values to a two-fold difference. This transformation, of course, would not be possible if some of the  $Y$  and  $X$

values are zero or negative. Though in this case, we can use  $\ln(Y_i + m)$  or  $\ln(X_i + m)$  where  $m$  is a positive value large enough to make all the values of  $Y$  or  $X$  positive.

In this question, our transformation is as follows:

$$\ln(PRICE_i) = \beta_0 + \beta_1 \ln(LOTSIZE_i) + \beta_2 \ln(SQRFT_i) + \beta_3 BDRMS_i + u_i$$

In R we can transform and obtain the BP and White tests as follows:

```
SQ2f_lm_log <- lm(log(PRICE) ~ log(LOTSIZE) + log(SQRFT) + BDRMS, data=property_df)
summary(SQ2f_lm_log)

# Breusch-Pagan test
bptest(SQ2f_lm_log)

# White test
bptest(SQ2f_lm_log, ~ log(LOTSIZE) + log(SQRFT) + BDRMS
      + I(log(LOTSIZE)^2) + I(log(SQRFT)^2) + I(BDRMS^2)
      + log(LOTSIZE)*log(SQRFT) + log(LOTSIZE)*BDRMS + log(SQRFT)*BDRMS, data=property_df)

# Special case of White test that conserves on degrees of freedom
bptest(SQ2f_lm_log, ~ fitted(SQ2f_lm_log) + poly(fitted(SQ2f_lm_log), 2))
```

Similarly, in STATA:

```
/* load the data */
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("heteroscedasticity") firstrow

/* generate the transformations */
generate lnPR = ln(PRICE)
generate lnLT = ln(LOTSIZE)
generate lnFT = ln(SQRFT)

/* run the regression */
regress lnPR lnLT lnFT BDRMS

/* run the Breusch-Pagan test */
hettest lnLT lnFT BDRMS, iid

/* run the White test */
imtest, white
```

Source		SS	df	MS	Number of obs	=	88
-----	+	-----	-----	-----	F(3, 84)	=	50.42
Model		5.1550402	3	1.71834673	Prob > F	=	0.0000
Residual		2.86256399	84	.034078143	R-squared	=	0.6430
-----	+	-----	-----	-----	Adj R-squared	=	0.6302
Total		8.0176042	87	.09215637	Root MSE	=	.1846
-----							
lnPR		Coefficient	Std. err.	t	P> t	[95% conf. interval]	
-----	+	-----	-----	-----	-----	-----	
lnLT		.1679666	.0382812	4.39	0.000	.0918403	.2440929

lnFT	.7002324	.0928653	7.54	0.000	.5155596	.8849051
BDRMS	.0369585	.0275313	1.34	0.183	-.0177905	.0917075
_cons	5.610714	.6512837	8.61	0.000	4.315565	6.905863

Breusch-Pagan/Cook-Weisberg test for heteroskedasticity

Assumption: i.i.d. error terms

Variables: lnLT lnFT BDRMS

H0: Constant variance

chi2(3) = 4.22  
Prob > chi2 = 0.2383

White's test

H0: Homoskedasticity

Ha: Unrestricted heteroskedasticity

chi2(9) = 9.55  
Prob > chi2 = 0.3882

Cameron & Trivedi's decomposition of IM-test

Source	chi2	df	p
Heteroskedasticity	9.55	9	0.3882
Skewness	1.69	3	0.6381
Kurtosis	2.57	1	0.1090
Total	13.81	13	0.3872

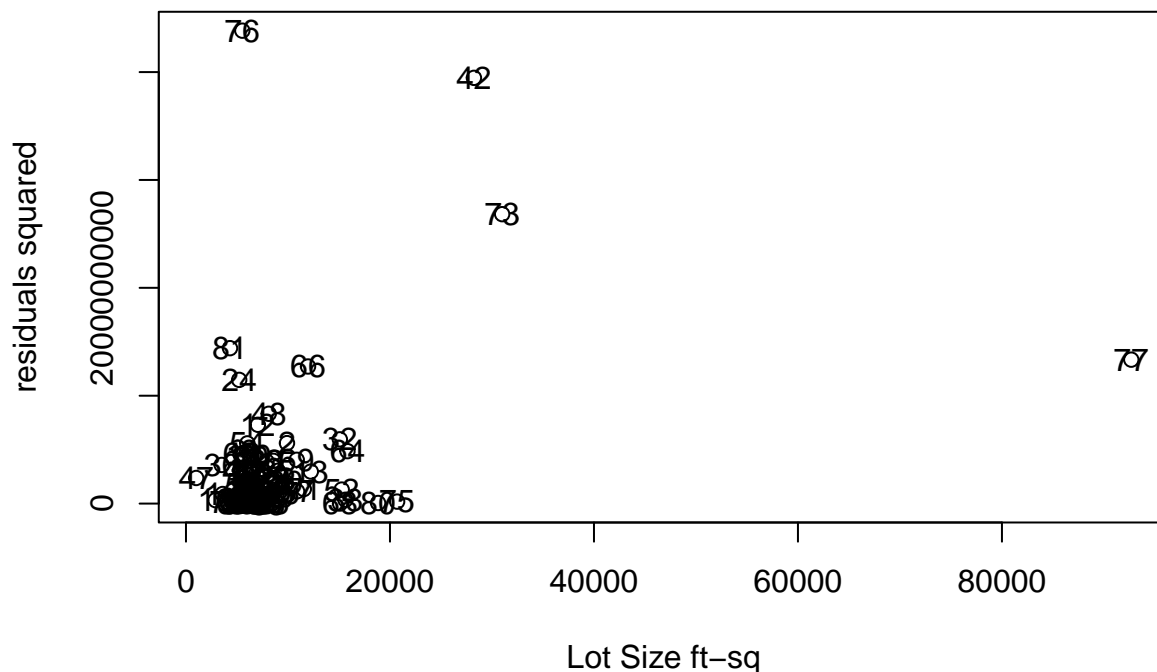
From the chi-squared values and their associated p-values we fail to reject the null hypothesis of homoskedasticity.

**Answer (ii):** Heteroskedasticity can also arise as a result of the presence of outliers. An outlier "is an observation from a different population to that generating the remaining sample observations."<sup>8</sup> Whether those outliers are included in the regression or not can substantially alter the results of the regression analysis. This is because OLS gives equal weight to every observation in the sample.

To see why these points were picked in the question first let's look at the scatter plot of *LOTSIZE*

```
plot(SQ2a_lm$residuals^2 ~ property_df$LOTSIZE,
     xlab="Lot Size ft-sq", ylab = "residuals squared")
text(SQ2a_lm$residuals^2 ~ property_df$LOTSIZE, label=property_df$Date)
```

<sup>8</sup>Gujarati and Porter (2009), Section 11.1



Here we can see that the 42nd, 73rd, 76th, and 77th data points are outliers.

We can now remove these and run the regression in R:

```
# create a new dataframe with outliers removed
property_df_nooutl <- property_df[-c(42,73,76,77),]

# run the regression using this new dataframe
SQ2f_lm_nooutl <- lm(PRICE ~ LOTSIZE + SQRFT + BDRMS, data = property_df_nooutl)
summary(SQ2f_lm_nooutl)

# run the BP test
bptest(SQ2f_lm_nooutl)

# run the White test
bptest(SQ2f_lm_nooutl, ~ BDRMS + LOTSIZE + SQRFT
      + I(BDRMS^2) + I(LOTSIZE^2) + I(SQRFT^2)
      + BDRMS*LOTSIZE + BDRMS*SQRFT + SQRFT*LOTSIZE,
      data = property_df_nooutl)

#Special case of White test that conserves on degrees of freedom
bptest(SQ2f_lm_nooutl, ~ fitted(SQ2f_lm_nooutl) + poly(fitted(SQ2f_lm_nooutl),2))
```

and in STATA

```
/* load the data */
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("heteroscedasticity") firstrow

/* run the regression removing the outliers */
regress PRICE LOTSIZE BDRMS SQRFT if Date != 42 & Date != 73 & Date != 76 & Date != 77

/* run the Breusch-Pagan test */
```

```

hettest LOTSIZE BDRMS SQRFT, iid
/* run the White test */
imtest, white

```

Source	SS	df	MS	Number of obs	=	84
				F(3, 80)	=	65.83
Model	3.7392e+11	3	1.2464e+11	Prob > F	=	0.0000
Residual	1.5146e+11	80	1.8933e+09	R-squared	=	0.7117
				Adj R-squared	=	0.7009
Total	5.2538e+11	83	6.3299e+09	Root MSE	=	43512

PRICE	Coefficient	Std. err.	t	P> t	[95% conf. interval]	
LOTSIZE	6.232396	1.529917	4.07	0.000	3.187764	9.277028
BDRMS	17387.21	6746.309	2.58	0.012	3961.626	30812.79
SQRFT	91.19175	11.08091	8.23	0.000	69.14004	113.2435
_cons	-8446.323	23975.92	-0.35	0.726	-56159.93	39267.29

Breusch-Pagan/Cook-Weisberg test for heteroskedasticity

Assumption: i.i.d. error terms

Variables: LOTSIZE BDRMS SQRFT

H0: Constant variance

```

chi2(3) = 5.97
Prob > chi2 = 0.1131

```

White's test

H0: Homoskedasticity

Ha: Unrestricted heteroskedasticity

```

chi2(9) = 9.26
Prob > chi2 = 0.4136

```

Cameron & Trivedi's decomposition of IM-test

Source	chi2	df	p
Heteroskedasticity	9.26	9	0.4136
Skewness	5.50	3	0.1386
Kurtosis	-226.48	1	1.0000
Total	-211.72	13	1.0000

From the chi-squared values and their associated p-values we fail to reject the null hypothesis of homoskedasticity.

(g) Considering your results from parts (e) and (f), which strategy for removing heteroskedasticity do you believe to be the best and why?

**Answer:** Table below summarizes the results of the BP and White tests for each strategy:

	F	Adj $R^2$	Heteroskedasticity Test		
				BP	White
Scaling by BDRMS	37.39	0.557	$\chi^2$	2.59	14.13
			$p$ -value	0.459	0.118
Scaling by LOTSIZE	480.77	0.943	$\chi^2$	5.96	8.74
			$p$ value	0.113	0.462
log transformation	50.42	0.63	$\chi^2$	4.22	9.55
			$p$ value	0.238	0.388
remove outliers in LOTSIZE	65.83	0.701	$\chi^2$	5.97	9.26
			$p$ value	0.113	0.414

Looking at the table, we can see that in all approaches we are able to tackle heteroskedasticity since in all of them we fail to reject the null hypothesis of homoskedasticity. However, the  $F$  value and the adjusted  $R^2$  is considerably higher for "scaling by *LOTSIZE*" strategy. This may be because the error variance is proportional to *LOTSIZE* and there are outliers in the *LOTSIZE* dimension. Scaling gets rid of both problems.

### QUESTION 3

(a) Briefly discuss the problem of autocorrelation - why might such a problem arise and what problems follow from the use of OLS?

**Answer:** We will briefly discuss the nature of the problem, why it occurs, and what problems arise from use of OLS with serially correlated errors.

What is the problem of autocorrelation?

Classical linear regression model assumes autocorrelation does not exist in the disturbances. That is, it assumes  $Cov(u_i, u_j | X_i, X_j) = 0$ ,  $i \neq j$ . What this means is that the disturbance term relating to any observation is not influenced by the disturbance term relating to any other observation. If there is such a dependence, then autocorrelation is present. That is,  $E(u_i, u_j) \neq 0$ .

Why autocorrelation occurs?

There are several reasons as to why autocorrelation occurs including inertia, specification bias, cobweb phenomenon, lags, 'manipulation' of data, data transformation, and nonstationarity.<sup>9</sup>

<sup>9</sup>Gujarati and Porter (2009), Section 12.1

inertia: Most econometrics time series has inertia or sluggishness in that the variables exhibit cyclical behavior where successive observations are likely to be interdependent.

specification bias: It may be the case that one of the variables of the true model is omitted or excluded from the model used. The omitted variable would then be part of the error or disturbance term, which would in turn reflect a systematic pattern.

cobweb phenomenon: The supply of many agricultural commodities reflect the cobweb phenomenon where supply reacts to price with a lag of one time period because supply decisions take time to implement. Thus at the beginning of this year's planting of crops, farmers are influenced by the price prevailing last year.

lags: In a time series regression, it is not uncommon to find that the dependent variable in the current variable also depends on the value of itself in the previous period. That is,  $Y_t = \beta_0 + \beta_1 X_t + \beta_2 Y_{t-1} + u_t$ . This is known as *autoregression*. If we omit the lagged term in this equation, the resulting error term will reflect a systematic pattern due to the influence of the lagged variable on the dependent variable.

'manipulation' of data: Smoothing out the data, say by averaging three monthly observations to obtain quarterly figures, may itself lend to a systematic pattern in the disturbances, thereby introducing autocorrelation. Similarly, interpolation or extrapolating from data for the missing values, or any data 'massaging' techniques may also introduce autocorrelation.

data transformation: Even if the level form  $Y_t = \beta_1 + \beta_2 X_t + u_t$  satisfies the OLS assumptions including no autocorrelation, in its first difference form  $\Delta Y_t = \beta_2 \Delta X_t + v_t$  the error term  $v_t = \Delta u_t$  is autocorrelated. Thus, certain transformations may induce autocorrelation.

nonstationarity: A time series is stationary if its characteristics such as mean, variance, and covariance are time invariant. If, however, they do change over time, then the time series are nonstationary, and the error term will exhibit autocorrelation.

#### Problems in estimating with autocorrelated errors:

Lets look at what happens in terms of unbiasedness, consistency, efficiency and goodness of fit.

Unbiasedness: Unbiasedness assumes nothing about the serial correlation of errors. As long as the explanatory variables are strictly exogenous then the  $\hat{\beta}_j$  are unbiased, regardless of the degree of serial correlation in the errors. This is analogous to the observation that heteroskedasticity in the errors does not cause bias in the  $\hat{\beta}_j$ .

consistency: Even when the data are weakly dependent, the  $\hat{\beta}_j$  are still consistent, though not necessarily unbiased if there is that dependency. Just like unbiasedness, this result does not depend on autocorrelation in the errors.

efficiency: Since the Gauss-Markov requires both homoskedasticity and serially uncorrelated errors, presence of autocorrelation in errors would make the OLS no longer BLUE. Even more importantly, the usual OLS standard errors and test statistics are not valid, *even asymptotically*.

#### Why not BLUE?

To see why OLS is no longer BLUE in the presence of serially correlated errors, consider the simple regression model  $Y_t = \beta_0 + \beta_1 X_t + u_t$  and assume that the error, or disturbance, terms are generated by the following mechanism:

$$u_t = \rho u_{t-1} + \varepsilon_t, \quad -1 < \rho < 1$$

where  $\rho$  is the *coefficient of autocovariance*, and  $\varepsilon_t$  are uncorrelated random variables that satisfies the OLS assumptions:

$$\mathbb{E}(\varepsilon_t) = 0 \quad ; \quad \text{Var}(\varepsilon_t) = \sigma^2 \quad ; \quad \text{Cov}(\varepsilon_t, \varepsilon_{t-s}) = 0 \quad , \quad s \neq 0.$$

Thus the error term in period  $t$  is equal to  $\rho$  times its value in the preceding period plus a white



noise error term. This scheme is known as a *Markov first-order autoregressive scheme*, or just a *first-order autoregressive scheme*, and usually denoted as **AR(1)**.

Also note that the population correlation coefficient between  $u_t$  and  $u_{t-1}$  is given by

$$\rho = \frac{\mathbb{E}\left([u_t - \mathbb{E}(u_t)][u_{t-1} - \mathbb{E}(u_{t-1})]\right)}{\sqrt{\text{Var}(u_t)}\sqrt{\text{Var}(u_{t-1})}} = \frac{\mathbb{E}(u_t u_{t-1})}{\text{Var}(u_{t-1})}$$

the latter equality is due to  $\mathbb{E}(u_t) = 0$  and  $\text{Var}(u_t) = \text{Var}(u_{t-1})$  because of the homoskedasticity assumption.  $\rho$  is also the slope coefficient in the regression of  $u_t$  on  $u_{t-1}$ .

Under the AR(1) scheme we have

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

$$\text{Cov}(u_t, u_{t-s}) = \mathbb{E}(u_t u_{t-s}) = \rho^s \frac{\sigma_\varepsilon^2}{1 - \rho^2}$$

$$\text{Corr}(u_t, u_{t-s}) = \rho^s$$

Since  $\rho$  is a constant between  $-1$  and  $1$ ,  $\text{Var}(u_t)$  is still homoskedastic. Notice, however,  $u_t$  is correlated not only with its immediate past value but its values several periods in the past.

If  $|\rho| = 1$  then the variances and covariances above are not defined. If  $|\rho| < 1$  then the AR(1) process is *stationary*, i.e. mean, variance, and covariance do not change over time. If  $|\rho| < 1$  it is also the case that the value of the covariance will decline as we go into distant past.

Considering the simple regression model  $Y_t = \beta_0 + \beta_1 X_t + u_t$  we know from previous supervisions that the OLS estimator of the slope coefficient is

$$\hat{\beta}_1 = \frac{\sum(X_t - \bar{X})(Y_t - \bar{Y})}{\sum(X_t - \bar{X})^2}$$

and its variance is

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum(X_t - \bar{X})^2}$$

Whereas under the AR(1) scheme, the variance of this estimator is

$$\begin{aligned} \text{Var}(\hat{\beta}_1)^{AR(1)} = & \frac{\sigma^2}{\sum(X_t - \bar{X})^2} \left[ 1 + 2\rho \frac{\sum(X_t - \bar{X})(X_{t-1} - \bar{X})}{\sum(X_t - \bar{X})^2} + 2\rho^2 \frac{\sum(X_t - \bar{X})(X_{t-2} - \bar{X})}{\sum(X_t - \bar{X})^2} \right. \\ & \left. + \dots + 2\rho^{n-1} \frac{\sum(X_t - \bar{X})(X_{t-n} - \bar{X})}{\sum(X_t - \bar{X})^2} \right] \end{aligned}$$

That is, the variance of  $\hat{\beta}_1$  under an AR(1) scheme is equal to its variance under the OLS times a term that depends on  $\rho$  as well as the sample autocorrelations between the values taken by the regressor  $X$  at various lags. In general, we cannot foretell whether  $\text{Var}(\hat{\beta}_1)$  is greater or less than  $\text{Var}(\hat{\beta}_1)^{AR(1)}$ .

If we assume that the regressor  $X$  also follows AR(1) scheme with a coefficient of autocorrelation, i.e. correlation between  $X_t$  and  $X_{t-1}$  given as

$$r = \frac{\sum(X_t - \bar{X})(X_{t-1} - \bar{X})}{\sum(X_t - \bar{X})^2}$$

then,  $\text{Var}(\hat{\beta}_1)^{AR(1)}$  reduces to

$$\text{Var}(\hat{\beta}_1)^{AR(1)} = \frac{\sigma^2}{\sum(X_t - \bar{X})^2} \left( \frac{1 + r\rho}{1 - r\rho} \right)$$

$$= \text{Var}(\hat{\beta}_1)^{OLS} \left( \frac{1 + r\rho}{1 - r\rho} \right)$$

If we continue to use the OLS estimator  $\hat{\beta}_1$  and adjust the usual variance formula by taking account of the AR(1) scheme,  $\hat{\beta}_1$  will still be linear and unbiased, but not efficient.

This finding is very similar to the finding that  $\hat{\beta}_1$  is less efficient in the presence of heteroskedasticity.

**(b) If  $Y_t = \alpha + \beta X_t + \varepsilon_t$  and  $\varepsilon_t = \rho\varepsilon_{t-1} + v_t$ , where  $\mathbb{E}(\varepsilon_t) = 0$ ,  $\rho \neq 0$ ,  $\text{Cov}(\varepsilon_t, v_t) = 0$  and  $v_t \sim i.i.d.N(0, \sigma^2)$ , show that:**

- i) if  $\varepsilon_t$  is homoskedastic,  $\text{Var}(\varepsilon_t) = \sigma^2/(1 - \rho^2)$ ;
- ii)  $\text{Cov}(\varepsilon_t, \varepsilon_{t-1}) \neq 0$ ;
- iii)  $\text{Corr}(\varepsilon_t, \varepsilon_{t-1}) = \rho$ , where *Corr* is the correlation coefficient.

**Answer (i):** Under AR(1) we have  $\varepsilon_t = \rho\varepsilon_{t-1} + v_t$  where  $v_t$  is a white noise error term. Since  $\varepsilon_t$  is homoskedastic, then

$$\mathbb{E}(\varepsilon_t) = \rho\mathbb{E}(\varepsilon_{t-1}) + \mathbb{E}(v_t) = 0$$

which means

$$\begin{aligned} \mathbb{E}(\varepsilon_t^2) &= \rho^2\mathbb{E}(\varepsilon_{t-1}^2) + \mathbb{E}(v_t^2) \\ \text{Var}(\varepsilon_t) &= \rho^2\text{Var}(\varepsilon_{t-1}) + \text{Var}(v_t) \end{aligned}$$

since  $\varepsilon$ s and  $v$ s are uncorrelated.

Because of homoskedasticity, notice that  $\text{Var}(u_t) = \text{Var}(u_{t-1}) = \sigma^2$ , and since  $v_t \sim i.i.d.N(0, \sigma^2)$  its variance is  $\text{Var}(v_t) = \sigma_v^2$ . Plugging these back into the expression above we get

$$\begin{aligned} \text{Var}(\varepsilon_t) &= \rho^2\text{Var}(\varepsilon_{t-1}) + \text{Var}(v_t) \\ &= \rho^2\text{Var}(\varepsilon_t) + \sigma_v^2 \\ \text{Var}(\varepsilon_t) - \rho^2\text{Var}(\varepsilon_t) &= \sigma_v^2 \\ \text{Var}(\varepsilon_t)(1 - \rho^2) &= \sigma_v^2 \\ \text{Var}(\varepsilon_t) &= \frac{\sigma_v^2}{1 - \rho^2}. \end{aligned}$$

as desired.

**Answer (ii):** To obtain the covariance, first we need to multiply the AR(1) scheme by  $\varepsilon_{t-1}$  and then we can take the expectations of the resulting expression, since that is the covariance. Accordingly,

$$\varepsilon_t = \rho\varepsilon_{t-1} + v_t$$

$$\varepsilon_t \varepsilon_{t-1} = \rho \varepsilon_{t-1}^2 + v_t \varepsilon_{t-1}$$

$$\mathbb{E}(\varepsilon_t \varepsilon_{t-1}) = \mathbb{E}(\rho \varepsilon_{t-1}^2 + v_t \varepsilon_{t-1})$$

$$\text{Cov}(\varepsilon_t, \varepsilon_{t-1}) = \mathbb{E}(\varepsilon_t \varepsilon_{t-1}) = \rho \mathbb{E}(\varepsilon_{t-1}^2) + \mathbb{E}(v_t \varepsilon_{t-1})$$

$$\text{Cov}(\varepsilon_t, \varepsilon_{t-1}) = \rho \mathbb{E}(\varepsilon_{t-1}^2) \quad \text{since } \text{Cov}(v_t, \varepsilon_{t-1}) = 0$$

$$\text{Cov}(\varepsilon_t, \varepsilon_{t-1}) = \rho \frac{\sigma_v^2}{1 - \rho^2} \quad \text{since } \text{Var}(\varepsilon_t) = \frac{\sigma_v^2}{1 - \rho^2}.$$

We can then continue in this fashion for the further past periods:

$$\text{Cov}(\varepsilon_t, \varepsilon_{t-2}) = \rho^2 \frac{\sigma_v^2}{1 - \rho^2}$$

$$\text{Cov}(\varepsilon_t, \varepsilon_{t-3}) = \rho^3 \frac{\sigma_v^2}{1 - \rho^2}$$

⋮

Since  $|\rho| \neq 0$ , then  $\text{Cov}(\varepsilon_t, \varepsilon_{t-1}) \neq 0$ .

**Answer (iii):** Since the correlation coefficient is the ratio of covariance to variance, we get

$$\text{Corr}(\varepsilon_t, \varepsilon_{t-1}) = \frac{\text{Cov}(\varepsilon_t, \varepsilon_{t-1})}{\text{Var}(\varepsilon_t)} = \frac{\rho \frac{\sigma_v^2}{1 - \rho^2}}{\frac{\sigma_v^2}{1 - \rho^2}} = \rho$$

and in future periods a similar approach would yield  $\text{Corr}(\varepsilon_t, \varepsilon_{t-2}) = \rho^2$ ,  $\text{Corr}(\varepsilon_t, \varepsilon_{t-3}) = \rho^3, \dots$

(c) Using the data from the worksheet ‘Autocorrelation’ estimate the following

$$gprice = \beta_0 + \beta_1 gwage_t + u_t \tag{7}$$

**Answer:** The question is asking for us to regress "growth of price" on "growth of wages". First, lets import the data and look at its summary.

In R:

```
wage_autocorr_df <- read_excel("../Data/sup4.xls", sheet = "autocorrelation")
str(wage_autocorr_df)
```

and in STATA via the command `describe`.

```
/* load the data and describe it*/
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow

describe
```

In R we see that the all the variables that start with *g*, i.e. all the growth columns, are coded as characters and in STATA the same are called strings. The data also seems to have missing values so we need to remove those as well.

In R we can do this with

```
# convert all characters to numeric
wage_autocorr_df <- wage_autocorr_df %>% mutate_if(is.character, as.numeric)

# remove N/As from `gprice` and `gwage` only
wage_autocorr_df <- wage_autocorr_df[!with(wage_autocorr_df,
                                           is.na(gprice) & is.na(gwage)),]

# look at the structure of the data again to confirm
str(wage_autocorr_df)
```

and we see that the number of observations for *gprice* and *gwage* now have reduced to 285 each.

The same can be done in STATA with

```
/* load the data*/
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow

/* change from string to numeric, and replace N/As */
quietly destring, replace
```

Now we can run the regression. In R

```
SQ3c_lm <- lm(gprice ~ gwage, data = wage_autocorr_df)
summary(SQ3c_lm)
```

Call:

```
lm(formula = gprice ~ gwage, data = wage_autocorr_df)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.0084151	-0.0021729	-0.0005211	0.0017191	0.0141002

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.0038241	0.0002745	13.93	< 0.0000000000000002 ***
gwage	0.1658072	0.0401518	4.13	0.0000479 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.003309 on 283 degrees of freedom

Multiple R-squared: 0.05683, Adjusted R-squared: 0.0535

F-statistic: 17.05 on 1 and 283 DF, p-value: 0.00004789

and in STATA

```
/* load the data*/
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow

/* change from string to numeric, and replace N/As */
quietly destring, replace

/* run the regression */
regress gprice gwage
```

(d) Test for serial correlation by (i) computing the Durbin-Watson  $d$  statistic, (ii) running a regression of  $\hat{u}_t$  on  $\hat{u}_{t-1}$  using residuals from part (c), (iii) using the Breusch-Godfrey test in Stata, (iv) using the alternative Durbin test in Stata (durbinalt). Discuss the relative advantages of each method.

**Answer (i):** A commonly used test for AR(1) autocorrelation is the Durbin-Watson  $d$  test<sup>10</sup> which is based on the OLS residuals.

#### Durbin-Watson $d$ test

The  $d$  statistic is the ratio of the sum of squared differences between successive residuals to the RSS:

$$DW = d = \frac{\sum_{i=2}^n (\hat{u}_i - \hat{u}_{i-1})^2}{\sum_{i=1}^n \hat{u}_i^2}$$

Note that the number of observations in the numerator of the  $d$  statistic is one less than the denominator, i.e.  $n - 1$ , because one observation is lost in taking successive differences.

#### Test procedure:

Unlike the  $t$ ,  $F$  or  $\chi^2$  tests, there is no unique critical value against which we would compare the  $d$  statistic for rejection decision. This is because, it is difficult to derive probability distribution of  $d$  statistic since it depends on  $\hat{u}$  which in turn depends on the  $X$ s.

To understand the decision process first expand the  $d$  statistic to get

$$d = \frac{\sum_{i=2}^n \hat{u}_i^2 + \sum_{i=2}^n \hat{u}_{i-1}^2 - 2 \sum_{i=2}^n \hat{u}_i \hat{u}_{i-1}}{\sum_{i=1}^n \hat{u}_i^2}$$

One thing to notice is that  $\sum \hat{u}_i^2$  and  $\sum \hat{u}_{i-1}^2$  differ in only one observation, so they are approximately equal. That is,  $\sum \hat{u}_i^2 \approx \sum \hat{u}_{i-1}^2$ . We can use this to rewrite the expression for  $d$  statistic:

$$d \approx 1 + 1 - 2 \left( \frac{\sum_{i=2}^n \hat{u}_i \hat{u}_{i-1}}{\sum_{i=1}^n \hat{u}_i^2} \right) = 2 \left( 1 - \frac{\sum_{i=2}^n \hat{u}_i \hat{u}_{i-1}}{\sum_{i=1}^n \hat{u}_i^2} \right)$$

<sup>10</sup>Durbin, J and Watson G S (1951) "Testing for Serial Correlation in Least-Squares Regression" *Biometrika* 38:159-177.

Also recall that in part (a) when discussing why OLS with autocorrelation is not BLUE, we noted that the definition of population correlation between  $u_t$  and  $u_{t-1}$  is

$$\rho = \frac{\mathbb{E}\left([u_t - \mathbb{E}(u_t)][u_{t-1} - \mathbb{E}(u_{t-1})]\right)}{\sqrt{\text{Var}(u_t)}\sqrt{\text{Var}(u_{t-1})}} = \frac{\mathbb{E}(u_t u_{t-1})}{\text{Var}(u_{t-1})}$$

the latter equality is due to  $\mathbb{E}(u_t) = 0$  and  $\text{Var}(u_t) = \text{Var}(u_{t-1})$  because of the homoskedasticity assumption.  $\rho$  is also the slope coefficient in the regression of  $u_t$  on  $u_{t-1}$ .

If we define  $\hat{\rho}$  as

$$\hat{\rho} = \frac{\sum_{i=2}^n \hat{u}_i \hat{u}_{i-1}}{\sum_{i=1}^n \hat{u}_i^2}$$

then the  $d$  statistic can be expressed as

$$d \approx 2(1 - \hat{\rho})$$

Since  $-1 \leq \hat{\rho} \leq 1$ , this means the bounds of  $d$  are

$$0 \leq d \approx 2(1 - \hat{\rho}) \leq 4.$$

That is, any estimated  $d$  statistic has to be between 0 and 4. We can use this for decision heuristics:

- If there are no autocorrelation in the residuals then  $\hat{\rho} = 0$  and  $d \approx 2$ . Therefore, if  $d = 2$  we can assume that there is no first-order autocorrelation;
- If there is perfect positive autocorrelation in the residuals then  $\hat{\rho} = 1$  and  $d \approx 0$ . Therefore, the closer  $d$  is to 0, the greater evidence of positive serial correlation.
- If there is perfect negative autocorrelation in the residuals then  $\hat{\rho} = -1$  and  $d \approx 4$ . Therefore, the closer  $d$  is to 4, the greater evidence of negative serial correlation.

Durbin-Watson then derived a lower bound  $d_L$  and an upper bound  $d_U$  within these limits of 0 and 4 such that if the computed  $d$  statistic lies outside of these bounds a decision can be made regarding the presence of negative or positive serial correlation as below:

If	Decision	Null Hypothesis
$0 < d < d_L$	Reject	No positive autocorrelation
$d_L \leq d \leq d_U$	No decision	No positive autocorrelation
$4 - d_L < d < 4$	Reject	No negative autocorrelation
$4 - d_U \leq d \leq 4 - d_L$	No decision	No negative autocorrelation
$d_U < d < 4 - d_U$	Do not reject	No autocorrelation

The upper and lower bounds depend on the number of observations,  $n$ , and the number of explanatory variables only. They do not depend on the values they take.

Accordingly, the steps for Durbin-Watson test is as follows:

1. Run the OLS and obtain the residuals;
2. Calculate  $d$ ;
3. For the given sample size  $n$  and the number of explanatory variables,  $k$ , obtain the critical values  $d_L$  and  $d_U$  from the manually published tables.

4. Follow the decision rules given in the table above.

Assumptions of  $d$  statistic:

1. The regression model includes the intercept term;
2. The explanatory variables, i.e. the  $X$ s, are nonstochastic, or fixed in repeated sampling;
3. The disturbances are generated by AR(1) scheme, i.e. it cannot be used to detect higher order autoregressive schemes;
4. The error term  $u_t$  is assumed to be normally distributed;
5. The regression model does not include lagged value(s) of the dependent variable as one of the explanatory variables, i.e. the test cannot be used for autoregressive models
6. There are no missing observations in the data.

Drawbacks of DW Test:

- It is only applicable when the above assumptions hold true;
- It cannot be used for testing for higher-autocorrelations than first-order;
- It contains zones of indecision where presence or absence of autocorrelation cannot be conclusively determined

We have already completed step 1 above in part (c) so we can now move to step 2 and compute the Durbin-Watson statistic.

For this, in R we use 'durbinWatsonTest()' function from the 'car' package, or 'dwtest' from the 'lmtest' package:

```
car::durbinWatsonTest(SQ3c_lm)
```

```
lag Autocorrelation D-W Statistic p-value
1      0.5940002      0.8032315      0
Alternative hypothesis: rho != 0
```

```
dwtest(SQ3c_lm)
```

Durbin-Watson test

```
data: SQ3c_lm
DW = 0.80323, p-value < 0.00000000000000022
alternative hypothesis: true autocorrelation is greater than 0
```

which gives us  $d = 0.8032315$  with a  $p$ -value of essentially 0.

Step 3 involves looking up the Durbin-Watson critical value table. It is difficult to find exact  $n = 285$  but we can approximate. For  $n = 280, k = 2$  we see that the critical values are  $d_L = 1.7969$  and  $d_U = 1.81123$ , and for  $n = 290, k = 2$  we see that the critical values are  $d_L = 1.80053$  and  $d_U = 1.81436$ . The values for  $n = 285, k = 2$  will be somewhere between these. In any case, Since  $d = 0.8032315$  is closer to 0 and lower

than either of the  $d_L$ , i.e.  $0 < d < d_L$ , we reject the null hypothesis and conclude that there seems to be a positive autocorrelation.

The same calculation in STATA can be done via the `estat dwatson` command

```
/* load the data*/
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow

/* change from string to numeric, and replace N/As */
quietly destring, replace

/* run the regression */
quietly regress gprice gwage

/* run Durbin-Watson test */
estat dwatson
```

**Answer (ii):** The second part of the question is asking us to regress the residuals to their one-lagged values.

In R:

```
SQ3d_ii_lm <- lm(SQ3c_lm$residuals ~ lag(SQ3c_lm$residuals, n=1))
summary(SQ3d_ii_lm)
```

Call:

```
lm(formula = SQ3c_lm$residuals ~ lag(SQ3c_lm$residuals, n = 1))
```

Residuals:

	Min	1Q	Median	3Q	Max
	-0.0125897	-0.0013615	-0.0001461	0.0015003	0.0157661

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.00001358	0.00015722	0.086	0.931
lag(SQ3c_lm\$residuals, n = 1)	0.59458883	0.04762707	12.484	<0.0000000000000002

(Intercept)

lag(SQ3c\_lm\$residuals, n = 1) \*\*\*

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.00265 on 282 degrees of freedom

(1 observation deleted due to missingness)

Multiple R-squared: 0.356, Adjusted R-squared: 0.3537

F-statistic: 155.9 on 1 and 282 DF, p-value: < 0.00000000000000022

Since  $d = 2(1 - \hat{\rho})$  and since  $\hat{\rho}$  is the slope of the regression coefficient of  $u_t$  on  $u_{t-1}$  - here 0.5945888 - we can calculate the  $d$  statistic as follows:



```
# rho is the slope coefficient
rho <- SQ3d_ii_lm$coefficients[2]

# calculate d
2*(1 - rho)
```

```
lag(SQ3c_lm$residuals, n = 1)
0.8108223
```

Here we get a very similar  $d$  value of 0.8108223 to 0.8032315 obtained in part (i), and would conclude the same way as we did in part (i) - reject the null hypothesis of no positive autocorrelation.

**Answer (iii):** To avoid some of the pitfalls of the Durbin-Watson  $d$  test of autocorrelation, a general test has been developed by Breusch-Godfrey.<sup>11</sup> The BG test can also be used if we think *gwage* is not strictly exogenous - a condition for the DW test.

#### Breusch-Godfrey (BG) test

Also known as Lagrange multiplier (LM) test, the BG test is a general test in the sense that it allows for

- nonstochastic regressors, such as the lagged values of the dependent variable;
- higher order autoregressive schemes such as AR(2), AR(3) etc; and
- simple or higher order *moving averages* of white noise errors terms; e.g. in regression  $Y_t = \beta_1 + \beta_2 X_t + u_t$ , the error term can be represented as  $u_t = \varepsilon_t + \eta_1 \varepsilon_{t-1} + \eta_2 \varepsilon_{t-2} + \dots + \eta_p \varepsilon_{t-p}$ , which represents a  $p$ -period moving average of the white noise error term  $\varepsilon_t$ .

#### Test procedure

Although the discussion can be extended to multiple regression, consider the simple regression model  $Y_t = \beta_0 + \beta_1 X_t + u_t$  with the error term following the  $p^{th}$ -order autoregressive, AR( $p$ ),<sup>a</sup> scheme:

$$u_t = \rho_1 u_{t-1} + \rho_2 u_{t-2} + \dots + \rho_p u_{t-p} + \varepsilon_t$$

where  $\varepsilon_t$  is the white noise error term. The null hypothesis is that there is no autocorrelation of any order:

$$\mathbb{H}_0 : \rho_1 = \rho_2 = \dots = \rho_p = 0$$

The BG or LM test is carried out via the following steps:

1. Run the OLS and obtain the residuals;
2. Regress the residuals,  $\hat{u}_t$ , on the original  $X_t$  and on its lagged values  $\hat{u}_{t-1}, \hat{u}_{t-2}, \dots, \hat{u}_{t-p}$ . If there are more than one  $X$  variable in the original model, then include them in this auxiliary regression also;
3. Obtain  $R^2$  from this auxiliary regression;
4. Since this auxiliary regression of the residuals have  $(n - p)$  observations given that  $p$  of them are used up in the model, Breusch and Godfrey have shown that  $R^2$  of this regression multiplied by the sample size  $(n - p)$  asymptotically follows the chi-square distribution

<sup>11</sup>Breusch, T S (1978) "Testing for Autocorrelation in Dynamic Linear Models", *Australian Economic Papers* 17:334-355 and Godfrey, L G (1978) "Testing Against General Autoregressive and Moving Average Error Models When the Regressor Includes Lagged Dependent Variables", *Econometrica* 46:1293-1302.

with  $p$  degrees of freedom. Therefore, if  $(n - p)R^2$  exceeds the critical  $\chi_p^2$  at the chosen level of significance, we reject the null hypothesis and conclude that at least one  $\rho$  is statistically significantly different from zero.

$$(n - p)R^2 \sim \chi_p^2$$

Drawbacks:

One drawback of the BG test is that the length of lag,  $p$ , cannot be specified a priori, which means some experimentation with  $p$  happens. Sometimes, Akaike and Schwarz information criteria can be used to select the lag length.

Another drawback is that the test assumes that the variance of disturbance is homoskedastic, i.e.  $\text{Var}(u_t) = \sigma^2$ .

---

<sup>a</sup>if  $p = 1$  then this test is called *Durbin's M test*.

To conduct the Breusch-Godfrey test for this question, we can follow these steps. step 1 has already been completed in part (i), so we will continue from step 2 onwards.

This can be done in R via 'bgtest()' function in 'lmtest' library:

```
bgtest(SQ3c_lm, order=1, data = wage_autocorr_df)
```

Breusch-Godfrey test for serial correlation of order up to 1

data: SQ3c\_lm

LM test = 105.4, df = 1, p-value < 0.00000000000000022

which gives us a  $\chi^2$  statistic of 105.4 with a  $p$  value of essentially 0. The critical value at  $\alpha = 0.05$  is  $\chi_1^2 = 3.84146$ . Since our test statistic exceeds that, we can reject the null hypothesis and conclude that  $\rho_1 \neq 0$  at 95% significance.

The same can be done in STATA using the `bgodfrey` command:

```
/* load the data*/
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow

/* change from string to numeric, and replace N/As */
quietly destring, replace

/* run the regression */
quietly regress gprice gwage

/* set variable `t` in the spreadsheet as time variable */
tsset t

/* run the Breusch-Godfrey test */
bgodfrey
```

**Answer (iii):** This final part of the question asks us to run the alternative Durbin test. The box below discusses the test briefly.

## Durbin's alternative test

Durbin's alternative test is in fact a Lagrange multiplier (LM) test but it is most easily computed with a Wald test on the coefficients of the lagged residuals in an auxiliary OLS regression. The auxiliary OLS regression regresses the residuals on their lags and all the explanatory variables in the original regression.

Consider the linear regression model

$$Y_t = \beta_1 X_{1t} + \cdots + \beta_k X_{kt} + u_t$$

where the covariates  $X_1$  through  $X_k$  are not assumed to be strictly exogenous<sup>a</sup> and  $u_t$  is assumed to be i.i.d. with finite variance. The process is also assumed to be stationary.

Test procedure

1. Run the OLS and obtain the residuals;
2. Run the auxiliary regression whereby the residuals  $\hat{u}_t$  are regressed on its lagged values and on the other regressors:

$$\hat{u}_t = \eta_1 \hat{u}_{t-1} + \cdots + \eta_p \hat{u}_{t-p} + \beta_1 X_{1t} + \cdots + \beta_k X_{kt} + \varepsilon_t$$

3. Durbin's alternative test is obtained by performing a Wald test on the null hypothesis that  $\mathbb{H}_0 : \eta_1 = \cdots = \eta_p = 0$

<sup>a</sup>Note that when there are only strictly exogenous regressors and  $p = 1$ , then this test is asymptotically equivalent to Durbin-Watson test.

In STATA we can run this test as follows:

```
/* load the data*/
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow

/* change from string to numeric, and replace N/As */
quietly destring, replace

/* run the regression */
quietly regress gprice gwage

/* set variable `t` in the spreadsheet as time variable */
tsset t

/* run Durbin's alternative test */
estat durbinalt
```

Time variable: t, 1 to 286

Delta: 1 unit

## Durbin's alternative test for autocorrelation

lags(p)	chi2	df	Prob > chi2
1	165.492	1	0.0000

H0: no serial correlation

which gives us a  $\chi^2$  value of 165.492.

The table below compares the results of each of the test:

Test	Statistic	Result
Durbin Watson	$d = 0.803$	Reject null
Manual calculation of DW	$d = 0.810$	Reject null
Breusch Godfrey	$\chi^2 = 105.39$	Reject null
Durbin's alternative	$\chi^2 = 165.49$	Reject null

Given that  $X$ s are not exogeneous, BG test is more appropriate here. However, it  $F$ -tests all slopes including all  $X$ s which is not really that interesting. So it may be preferable to use Durbin's alternative test instead, which is almost identical to the BG test but it does not test whether all the slope coefficients are zero, only whether the coefficients of the lagged residuals are zero.

(e) Estimate the model using the Cochrane-Orcutt (two step) method. Check that these results are the same as those obtained by using the `prais` command in Stata (with the `corc` and `twostep` options). Does any autocorrelation remain?

**Answer:** This question is looking at the ways in which we can deal with autocorrelation, especially the lack of efficiency of the OLS estimators. Consider a simple regression model  $Y_t = \beta_0 + \beta_1 X_t + u_t$  with the errors following the AR(1) process

$$u_t = \rho u_{t-1} + \varepsilon_t$$

Since there is autocorrelation in the disturbance term, we can transform the equation to try to eliminate the problem. We can do this in such a way that the error term of our transformed equation is the white noise  $\varepsilon_t$ . To obtain that, the error term of the transformed equation then has to equal to  $u_t - \rho u_{t-1}$ . Therefore, for  $t \geq 2$  we can lag the equation by one period, multiply it by  $\rho$ , and subtract it from the original equation.

$$\begin{aligned}\rho Y_{t-1} &= \beta_0 \rho + \beta_1 \rho X_{t-1} + \rho u_{t-1} \\ Y_t - \rho Y_{t-1} &= \beta_0 - \beta_0 \rho + \beta_1 X_t - \beta_1 \rho X_{t-1} + u_t - \rho u_{t-1} \\ Y_t - \rho Y_{t-1} &= \beta_0(1 - \rho) + \beta_1(X_t - \rho X_{t-1}) + \varepsilon_t \\ Y_t^* &= \beta_0^* + \beta_1 X_t^* + \varepsilon_t\end{aligned}$$

where  $Y_t^* = Y_t - \rho Y_{t-1}$ ,  $X_t^* = X_t - \rho X_{t-1}$ ,  $\beta_0^* = \beta_0(1 - \rho)$ , and  $-1 < \rho < 1$ . With the white noise error, the model is now free from autocorrelation, as desired.

With multiple regressors with  $u_t$  following an AR(1) process, the transformed model can be expressed as

$$Y_t^* = \beta_0^* + \beta_1 X_{1t}^* + \beta_2 X_{2t}^* + \cdots + \beta_k X_{kt}^* + \varepsilon_t.$$

Since now the error term satisfies the Gauss-Markov assumptions which means that if we knew  $\rho$ , we could estimate  $\beta_0^*$  and  $\beta_1$  by regressing  $Y_t^*$  on  $X_t^*$ . In effect, running  $Y_t^* = \beta_0^* + \beta_1 X_t^* + \varepsilon_t$  is tantamount to using generalized least squares (GLS) discussed in the Supplementary Questions 1(b) above, where it was posited that GLS is the application of OLS to a transformed model that satisfies the Gauss-Markov assumptions.

This equation,  $Y_t^* = \beta_0^* + \beta_1 X_t^* + \varepsilon_t$ , is known as *generalized difference equation*, or *quasi-differenced equation*.

### Prais-Winsten transformation

Notice that in the differencing procedure above we lose one observation because the lagged variables are not defined for the first observation. This loss of observation may not be too serious in large samples but can make considerable difference in the results of small samples. However, the OLS estimators from  $Y_t^* = \beta_0^* + \beta_1 X_t^* + \varepsilon_t$  is not quite BLUE because the first period is not used.

A quick and straight-forward fix for this is to write the equation for  $t = 1$  as

$$Y_1 = \beta_0 + \beta_1 X_1 + u_1.$$

This is a fix because the error term of the first observation,  $u_1$ , is uncorrelated with the white noise error  $\varepsilon_t$ . Therefore we can add this equation to  $Y_t^* = \beta_0^* + \beta_1 X_t^* + \varepsilon_t$ .

However, this introduces a problem. With the inclusion of the first observation, the error variances will not be homoskedastic. This is because the variance of  $u_t$  is

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

as derived in Supplementary Question 3(b). Since  $|\rho| < 1$ ,

$$Var(u_t) = \frac{\sigma_\varepsilon^2}{1 - \rho^2} > \sigma_\varepsilon^2 = Var(\varepsilon_t)$$

In order to equate the two variances therefore, we need to adjust the equation of the first observation by  $\sqrt{1 - \rho^2}$  whereby

$$\sqrt{1 - \rho^2} Y_1 = \beta_0 \sqrt{1 - \rho^2} + \beta_1 \sqrt{1 - \rho^2} X_1 + \sqrt{1 - \rho^2} u_1$$

so that the error variance of this equation becomes

$$Var(\sqrt{1 - \rho^2} u_1) = (1 - \rho^2) Var(u_t) = (1 - \rho^2) \frac{\sigma_\varepsilon^2}{1 - \rho^2} = \sigma_\varepsilon^2 = Var(\varepsilon_t).$$

Adding this adjusted equation for first observation to  $Y_t^* = \beta_0^* + \beta_1 X_t^* + \varepsilon_t$  then yields OLS estimators that are BLUE.

One of the points that was made above was that if we knew  $\rho$ , we could estimate  $\beta_0^*$  and  $\beta_1$  by regressing  $Y_t^*$  on  $X_t^*$ . Cochrane-Orcutt method is a way to estimate  $\rho$ .

### Cochrane-Orcutt Two-Step Procedure

Early regression applications tended to use the Cochrane-Orcutt iterative process to provide an estimate for  $\rho$ . One advantage of this process is that it can be used to estimate not only AR(1) scheme but higher order autoregressive schemes such as  $\hat{u}_t = \hat{\rho}_1 \hat{u}_{t-1} + \hat{\rho}_2 \hat{u}_{t-2} + v_t$ , which is AR(2).

The procedure for this method is as follows:

1. Run the OLS on the untransformed model;
2. Regress  $\hat{u}_t = \rho \hat{u}_{t-1} + \varepsilon_t$  to obtain an estimate of  $\rho$ ;
3. Use this estimation,  $\hat{\rho}$ , to fit  $Y_t^* = \beta_0^* + \beta_1 X_t^* + \varepsilon_t$ .

By using two linear regressions, we can obtain a BLUE estimate.

For this question, therefore we need to transform the following model:

$$gprice_t^* = \beta_0^* + \beta_1 gwage_t^* + \varepsilon_t$$

$$gprice_t - \hat{\rho} gprice_{t-1} = \beta_0(1 - \hat{\rho}) + \beta_2(gwage_t - \hat{\rho} gwage_{t-1}) + u_t - \hat{\rho} u_{t-1}$$

In supplementary question 3(d)(ii) above we already obtained an estimate for  $\rho$  as 0.5945888. We can use this to transform the model and run the regression.

In R:

```
# transform the data
wage_autocorr_df <- wage_autocorr_df %>%
  mutate(gprice_star = gprice - rho*lag(gprice, 1),
         gwage_star = gwage - rho*lag(gwage,1),
         beta0_star = 1-rho)

# next we need to adjust this for the first observation
wage_autocorr_df$gprice_star[wage_autocorr_df$t == 2] <-
  sqrt((1- rho^2))*wage_autocorr_df$gprice[wage_autocorr_df$t == 2]

wage_autocorr_df$gwage_star[wage_autocorr_df$t == 2] <-
  sqrt((1- rho^2))*wage_autocorr_df$gwage[wage_autocorr_df$t == 2]

wage_autocorr_df$beta0_star[wage_autocorr_df$t == 2] <-
  sqrt((1- rho^2))

# run the regression on the transformed variables
SQ3e_lm <- lm(gprice_star ~ gwage_star + beta0_star + 0, data = wage_autocorr_df)
summary(SQ3e_lm)

# or alternatively use `prais_winsten()` formula from `prais` package
library(prais)
prais_winsten(gprice ~ gwage, data = wage_autocorr_df, index = "t", twostep = TRUE)
```

and in STATA:

```
/* load the data*/
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow

/* change from string to numeric, and replace N/As */
quietly destring, replace

/* transform the variables */
quietly generate gprice_star = gprice - 0.5945888 * gprice_1
quietly generate gwage_star = gwage - 0.5945888 * gwage_1
quietly generate beta0_star = 1 - 0.5945888

/* adjust the first observation */
quietly replace gprice_star=((1 - 0.5945888^2)^(1/2))*gprice if t==2
quietly replace gwage_star=((1 - 0.5945888^2)^(1/2))*gwage if t==2
quietly replace beta0_star=((1 - 0.5945888^2)^(1/2)) if t==2
```

```
/* run the regression */
regress gprice_star gwage_star beta0_star, noc
```

But this can instead be done automatically with the `prais` command in STATA:

```
/* load the data*/
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow

/* change from string to numeric, and replace N/As */
quietly destring, replace

/* set the time index with volumn "t" */
tsset t

/* run the `prais` command */
prais gprice gwage, twostep
```

Time variable: t, 1 to 286  
Delta: 1 unit

Iteration 0: rho = 0.0000  
Iteration 1: rho = 0.5946

Prais-Winsten AR(1) regression with twostep estimates

Source	SS	df	MS	Number of obs	=	285
Model	3.0586e-06	1	3.0586e-06	F(1, 283)	=	0.47
Residual	.00183807	283	6.4949e-06	Prob > F	=	0.4931
				R-squared	=	0.0017
				Adj R-squared	=	-0.0019
Total	.001841128	284	6.4828e-06	Root MSE	=	.00255

	gprice	Coefficient	Std. err.	t	P> t	[95% conf. interval]
	gwage	.0405599	.0254657	1.59	0.112	-.0095663 .0906861
	_cons	.0043858	.0003898	11.25	0.000	.0036186 .005153
	rho	.5945966				

Durbin-Watson statistic (original) = 0.803231  
Durbin-Watson statistic (transformed) = 2.182960

and to obtain the Cochrane-Orcutt Two-Step Procedure:

```
/* load the data*/
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow

/* change from string to numeric, and replace N/As */
```

```
quietly destring, replace

/* set the time index with volumn "t" */
tsset t

/* run the `prais` command */
prais gprice gwage, corc twostep
```

Time variable: t, 1 to 286  
Delta: 1 unit

Iteration 0: rho = 0.0000  
Iteration 1: rho = 0.5946

Cochrane-Orcutt AR(1) regression with twostep estimates

Source	SS	df	MS	Number of obs	=	284
				F(1, 282)	=	2.34
Model	.000015099	1	.000015099	Prob > F	=	0.1271
Residual	.001818457	282	6.4484e-06	R-squared	=	0.0082
				Adj R-squared	=	0.0047
Total	.001833555	283	6.4790e-06	Root MSE	=	.00254

gprice	Coefficient	Std. err.	t	P> t	[95% conf. interval]
gwage	.0388561	.0253931	1.53	0.127	-.011128 .0888403
_cons	.0044698	.0003913	11.42	0.000	.0036994 .0052401
rho	.5945966				

Durbin-Watson statistic (original) = 0.803231  
Durbin-Watson statistic (transformed) = 2.194615

Thus, we obtain  $\beta_1 = 0.0388568$  and  $\beta_0(1 - 0.5945888) = 0.0018121$ , which means  $\beta_0 = 0.00446978$ .

Notice that the transformed DW  $d$  statistic is 2.194615. The critical values are  $d_L = 2.1887$  and  $d_U = 2.203$  which means this falls within the “No decision” rule.

To see if we reduced autocorrelation we can first regress the residuals from the transformed model's regression against their 1-lagged values, i.e.  $\varepsilon_t = \gamma_0 + \gamma_1 \varepsilon_{t-1} + v_t$ :

```
summary(lm(SQ3e_lm$residuals ~ lag(SQ3e_lm$residuals, 1)))
```

Error in eval(predvars, data, env): object 'SQ3e\_lm' not found

The estimate of the coefficient of the lagged value,  $\gamma_1$ , does not seem to be significant since it is significant only at  $\alpha = 0.05$ . We can check if this is also the case for higher-order lags, i.e.  $\varepsilon_t = \gamma_0 + \gamma_1 \varepsilon_{t-1} + \gamma_2 \varepsilon_{t-2} + \gamma_3 \varepsilon_{t-3} + v_t$ :

```
summary(lm(SQ3e_lm$residuals ~ lag(SQ3e_lm$residuals, 1)
+ lag(SQ3e_lm$residuals, 2) + lag(SQ3e_lm$residuals, 3)))
```



```
Error in eval(predvars, data, env): object 'SQ3e_lm' not found
```

We see that the second and third lags,  $\gamma_2$  and  $\gamma_3$ , are significant at  $\alpha = 0.01$ . which suggests there may be AR(2) and AR(3) process in the error terms. We can also check this via Breusch-Godfrey test:

```
bgtest(SQ3e_lm, order = 1, data = property_autocorr_df)
```

```
Error in eval(expr, envir, enclos): object 'SQ3e_lm' not found
```

```
bgtest(SQ3e_lm, order = 2, data = property_autocorr_df)
```

```
Error in eval(expr, envir, enclos): object 'SQ3e_lm' not found
```

```
bgtest(SQ3e_lm, order = 3, data = property_autocorr_df)
```

```
Error in eval(expr, envir, enclos): object 'SQ3e_lm' not found
```

which gives a  $\chi^2$  values of 8.8743 and 16.137 for lags 2 and 3, with  $p$  values suggesting significance at least at  $\alpha = 0.05$  level.

So based on the above tests it is possible that some autocorrelation may still be present.

The STATA equivalent of above commands are:

```
/* load the data*/
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow

/* change from string to numeric, and replace N/As */
quietly destring, replace

/* transform the variables */
generate gprice_star = gprice - 0.5945888 * gprice_1
generate gwage_star = gwage - 0.5945888 * gwage_1

/* run the regression */
quietly regress gprice_star gwage_star

/* regress residuals against its lags */
tsset t
predict U, residuals
regress U L.U L2.U L3.U

/* Run Breusch-Godfrey manually */
display e(N)*e(r2)

/* run regression with residuals and regressors */
regress gprice_star gwage_star L.U L2.U L3.U

/* Run Breusch-Godfrey */
bgodfrey, lags (1 2 3)
```

(f) Use the 'ac' command to graph the autocorrelations of the residuals from equation (7) on page 43 that you used in part (c), and comment on the relevance of this graph for the results above.

```
/* load the data*/
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow

/* change from string to numeric, and replace N/As */
quietly destring, replace

/* transform the variables */
generate gprice_star = gprice - 0.5945888 * gprice_1
generate gwage_star = gwage - 0.5945888 * gwage_1

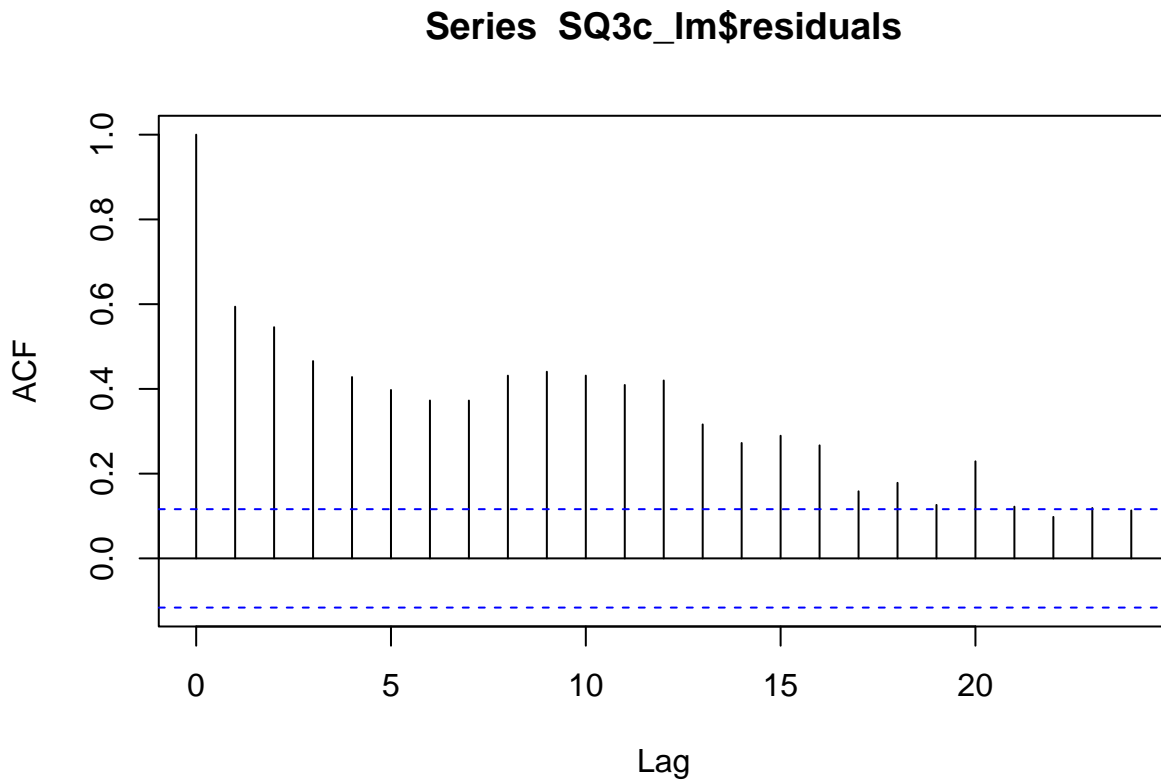
/* run the regression */
quietly regress gprice_star gwage_star

/* regress residuals against its lags */
tsset t
predict U, residuals

/* generate autocorrelations and obtain the graph */
quietly corrgram U, lags(40)
ac U, lags(40)
```

and in R we use the `Acf()` function from the `forecast` package:

```
library(forecast)
acf(SQ3c_lm$residuals)
```



Which shows that there is a lot of autocorrelation present in the residuals of equation 7, and much of it is of a higher order than the ones we have attempted to deal with so far.

Note that the dashed lines are the Bartlett's standard errors for MA(q) 95% confidence intervals.

(g) Now estimate the following distributed lag model:

$$gprice_t = \beta_0 + \beta_1 gwage_t + \beta_2 gwage_{t-1} + \cdots + \beta_{13} gwage_{t-12} + \varepsilon_t \quad (8)$$

Comment on your results. Is there autocorrelation in this model?

**Answer:** The spreadsheet has these lagged values as columns already. So we just need to regress them.

In R:

```
SQ3g_lm <- lm(gprice ~ gwage + gwage_1 + gwage_2 + gwage_3 + gwage_4 + gwage_5
              + gwage_6 + gwage_7 + gwage_8 + gwage_9 + gwage_10 + gwage_11
              + gwage_12, data = property_autocorr_df)
```

Error in eval(mf, parent.frame()): object 'property\_autocorr\_df' not found

```
Acf(SQ3g_lm$residuals)
```

```
Error in eval(expr, envir, enclos): object 'SQ3g_lm' not found
```

We can see that the autocorrelation is reduced but not eliminated, especially in lower lags.

As before, we can also use Breusch Godfrey to check for autocorrelation:

```
bgtest(SQ3g_lm)
```

```
Error in eval(expr, envir, enclos): object 'SQ3g_lm' not found
```

with  $\chi^2 = 70.241$  and a  $p$  value that is practically zero, we would reject the null hypothesis that there is no autocorrelation in the residuals.

---

(h) Run model (8) on page 59 again, this time including the residuals from this equation as lagged values. Experiment with the number of lagged terms you include. How many lagged terms are required to make the autocorrelation disappear?

```
SQ3h_lm <- update(SQ3g_lm, . ~ . + SQ3g_lm$residuals + lag(SQ3g_lm$residuals, 1))
Acf(SQ3h_lm$residuals)
```

```
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow

quietly deststring, replace

tsset t

quietly regress gprice gwage gwage_1 gwage_2 gwage_3 gwage_4 gwage_5 gwage_6 gwage_7 gwage_8 gwage_9 gwage_10
predict E, residuals

/* run the regression again with its residuals and their lags */

regress gprice gwage gwage_1 gwage_2 gwage_3 gwage_4 gwage_5 gwage_6 gwage_7 gwage_8 gwage_9 gwage_10

/* Breusch Godfrey */
bgodfrey, lags (1 2 3 4)

/* Durbin's alternative test */
estat durbinalt, lags (1 2 3 4)
```

Time variable: t, 1 to 286  
Delta: 1 unit

(13 missing values generated)

Source	SS	df	MS	Number of obs	=	271
				F(15, 255)	=	17.87
Model	.00156737	15	.000104491	Prob > F	=	0.0000
Residual	.001490823	255	5.8464e-06	R-squared	=	0.5125
				Adj R-squared	=	0.4838
Total	.003058193	270	.000011327	Root MSE	=	.00242

gprice	Coefficient	Std. err.	t	P> t	[95% conf. interval]	
gwage	.0741411	.0445721	1.66	0.097	-.0136353	.1619174
gwage_1	.0753577	.0334012	2.26	0.025	.0095805	.141135
gwage_2	.0333952	.0333378	1.00	0.317	-.0322572	.0990476
gwage_3	.0393801	.0332849	1.18	0.238	-.0261682	.1049285
gwage_4	.0821246	.0334291	2.46	0.015	.0162924	.1479569
gwage_5	.1148704	.0334868	3.43	0.001	.0489244	.1808164
gwage_6	.0950559	.0333798	2.85	0.005	.0293206	.1607912
gwage_7	.0977234	.0333518	2.93	0.004	.0320432	.1634035
gwage_8	.1051133	.0334565	3.14	0.002	.0392271	.1709994
gwage_9	.1617887	.0333054	4.86	0.000	.0962001	.2273774
gwage_10	.1139434	.0332565	3.43	0.001	.048451	.1794358
gwage_11	.1038893	.0334864	3.10	0.002	.0379442	.1698345
gwage_12	.0428162	.0443648	0.97	0.335	-.0445518	.1301843
E						
L1.	.4038747	.0617549	6.54	0.000	.2822601	.5254892
L2.	.2082131	.0614613	3.39	0.001	.0871766	.3292496
_cons	-.0007536	.0004821	-1.56	0.119	-.0017031	.0001958

Breusch-Godfrey LM test for autocorrelation

lags(p)	chi2	df	Prob > chi2
1	1.376	1	0.2408
2	1.452	2	0.4839
3	2.655	3	0.4479
4	3.080	4	0.5445

H0: no serial correlation

Durbin's alternative test for autocorrelation

lags(p)	chi2	df	Prob > chi2
---------	------	----	-------------

1		1.296	1	0.2549
2		1.363	2	0.5060
3		2.493	3	0.4765
4		2.886	4	0.5771

---

H0: no serial correlation

Notice that the  $\chi^2$  statistics are getting larger for higher-order lags but this is mainly because of degrees of freedom being used up and not because they are becoming less significant (notice the  $p$ -values are increasing too.)

Overall, we used residuals and its one lag as well for the model.

---

(i) Using this last specification estimate the long-run elasticity between prices and wages and test that it is unity.

```
quietly cd ..
quietly import excel using Data/sup4.xls, sheet("autocorrelation") firstrow
quietly destring, replace
tsset t
quietly regress gprice gwage gwage_1 gwage_2 gwage_3 gwage_4 gwage_5 gwage_6 gwage_7 gwage_8 gwage_9
predict E, residuals

/* run the regression again with its residuals and their lags */
quietly regress gprice gwage gwage_1 gwage_2 gwage_3 gwage_4 gwage_5 gwage_6 gwage_7 gwage_8 gwage_9

/* obtain the linear combination of parameter estimates */
lincom gwage + gwage_1 + gwage_2 + gwage_3 + gwage_4 + gwage_5 + gwage_6 + gwage_7 + gwage_8 + gwage_9

/* carry out a Wald test that linear combination equals to 1 */
test gwage + gwage_1 + gwage_2 + gwage_3 + gwage_4 + gwage_5 + gwage_6 + gwage_7 + gwage_8 + gwage_9
```

Time variable: t, 1 to 286  
Delta: 1 unit

(13 missing values generated)

```
( 1)  gwage + gwage_1 + gwage_2 + gwage_3 + gwage_4 + gwage_5 + gwage_6 + gwage_7
      + gwage_8 + gwage_9 + gwage_10 + gwage_11 + gwage_12 = 0
```

---

gprice	Coefficient	Std. err.	t	P> t	[95% conf. interval]
--------	-------------	-----------	---	------	----------------------

(1)	1.139599	.0940975	12.11	0.000	.9542923	1.324907
-----	----------	----------	-------	-------	----------	----------

$$(1) \quad \text{gwage} + \text{gwage}_1 + \text{gwage}_2 + \text{gwage}_3 + \text{gwage}_4 + \text{gwage}_5 + \text{gwage}_6 + \text{gwage}_7 + \text{gwage}_8 + \text{gwage}_9 + \text{gwage}_{10} + \text{gwage}_{11} + \text{gwage}_{12} = 1$$

F( 1, 255) =	2.20
Prob > F =	0.1392

or alternatively we can carry out a  $t$ -test on the coefficient which should be the square root of the Wald's F statistic above:

$$(1.139599-1)/0.0940975$$

[1] 1.483557

$$((1.139599-1)/0.0940975)^2$$

[1] 2.200941

In either approach, we cannot reject the null hypothesis that the long run elasticity between prices and wages is unity.