

ECONOMETRICS 1 / APPLIED PROBLEM SET 1

Topic: Introduction to Multivariate Regression

- This class is an introduction to the **menu-driven features** of EViews 6. The data is available on the website as `cigarette.txt`.
- The problem set covers (i) file import and workfile save, (ii) statistical analysis, (iii) creation of groups, (iv) ordinary least squares, (v) (robust) inference, and (vi) diagnostic testing.
- We use data on (i) the log of cigarette consumption (in packs) per person of smoking age ( $> 16$  years) for 46 U.S. states in 1992: LNC ( $\ln(C)$ ), (ii) the log real price of cigarettes in each state, normalized at 1983\$ per pack: LNP ( $\ln(P)$ ), and (iii) the log of real disposable income per capita in each state, in 1983\$1000: LNY ( $\ln(Y)$ ).
- Perform all of the steps described in Figures 1-85 (answering any corresponding questions, and carefully considering the methods and output, noting any new or confusing tools), and end by responding in full to the following questions:

1. Regress log consumption on log prices (eq01):

$$\widehat{\ln(C_i)} \approx [5.09] + [-1.20] \ln(P_i),$$

$$\begin{aligned} y_i &= x_i' \beta + u_i \\ \hat{y}_i &= x_i' \hat{\beta} \\ \hat{u}_i &= y_i - \hat{y}_i \end{aligned}$$

with estimated residuals  $\hat{u}_i^{(1)} \approx \ln(C_i) + [-5.09] + [1.20] \ln(P_i)$ .

What is the estimated price elasticity of consumption? Answer:

$$\frac{\% \Delta \ln C}{\% \Delta \ln P} \approx \frac{100}{100} \frac{\partial \widehat{\ln(C)} / C}{\partial \ln(P) / P} = \frac{\partial \widehat{\ln(C)}}{\partial \ln(P)} \approx [-1.20] \quad ; \quad \frac{\partial \ln C}{\partial C} = \frac{1}{C}$$

and so a [ 10% ] ~~increase/decrease~~ in price results in a [ 12% ] ~~in-~~  
~~crease/decrease~~ in consumption.

2. Regress log consumption on log prices and log income (eq02): no "satiation"  
ie.  $Y \uparrow \Rightarrow c \uparrow$

$$\widehat{\ln(C_i)} \approx [ 4.30 ] + [ -1.34 ] \ln(P_i) + [ 0.17 ] \ln(Y_i).$$

3. Regress log income on log prices (eq03):

$P \uparrow 10\%$   $c \downarrow 13\%$   
 $Y \uparrow 10\%$   $c \uparrow 2\%$  (insignificant)

$$\widehat{\ln(Y_i)} \approx [ 4.61 ] + [ 0.81 ] \ln(P_i),$$

$Y \perp P$   $\rightarrow$   
with estimated residuals  $\widehat{u}_i^{(3)} \approx \ln(Y_i) + [ -4.61 ] + [ -0.81 ] \ln(P_i)$ .

4. Regress log consumption on  $\widehat{u}_i^{(3)}$  (eq04):

$$\widehat{\ln(C_i)} \approx [ 4.85 ] + [ 0.17 ] \widehat{u}_i^{(3)}$$

5. Regress  $\widehat{u}_i^{(1)}$  on  $\widehat{u}_i^{(3)}$  (eq05):

$$\widehat{u}_i^{(1)} \approx [ 1.2 \times 10^{-15} ] + [ 0.17 ] \widehat{u}_i^{(3)}.$$

6. Regress log consumption on log income, log income squared, and log prices (eq06):

$$\widehat{\ln(C_i)} \approx [ -34.77 ] + [ 16.43 ] \ln(Y_i) + [ -1.69 ] (\ln(Y_i))^2 + [ -1.35 ] \ln(P_i),$$

↓ (rounding)

so that a [ 10% ] ~~increase/decrease~~ in price results in a [ 13% ] ~~in-~~

crease/decrease in consumption. The income elasticity of consumption is now:

$$\frac{\partial \ln(\hat{C}_i)}{\partial \ln(Y_i)} \approx [16.43] + [-3.38] \ln(Y_i),$$

see histogram of  
ln Y. very few  
observations are  
< 4.57

which is greater than (equal to) less than 1 as  $\ln(Y_i)$  is less than (equal to) greater than  $[4.57]$  (as  $Y_i$  is less than (equal to) greater than  $[\$96,147]$ ): use the full available accuracy on the estimated coefficients when performing this computation. Interpret your findings carefully.

$\Rightarrow \frac{\partial \ln C}{\partial \ln Y} \leq 1$

eq02  $C \sim P, Y$

eq04  $C \sim Y \perp P$

eq05  $C \perp P \sim Y \perp P$

7. Compare equations eq02, eq04 and eq05. In eq02,  $[0.17]$  quantifies the impact of log income on log consumption. In eq04,  $\hat{u}_i^{(3)}$  is the part of log income not explained by log price, and so  $[0.17]$  quantifies the impact on log consumption of that part of log income not explained by log price. In eq05,  $\hat{u}_i^{(1)}$  is the part of log consumption not explained by log price, and so  $[0.17]$  quantifies the impact on that part of log consumption not explained by log price, of that part of log income not explained by log price. Carefully explain the intuition behind these results.

— each model "controls for" P in a different way  $\Rightarrow$  same result for Y coefficient

- Special attention should be paid to observations 3 (Arkansas), 15 (Kentucky) and 40 (Utah): Arkansas and Kentucky have particularly high sales, Kentucky is a producer with rather low prices, and Utah has especially low sales because of its high Mormon population (which bans smoking). It is important to build a deep understanding of the structure and peculiarities of your data before you start modelling.

data  $i=1, 2, \dots, n$   
 $x_1, x_2, \dots, x_n$   
 ordered data

$x_{(1)} \leq \dots \leq x_{(n)}$

minimum

$x_{(1)}$

maximum

$x_{(n)}$

mean

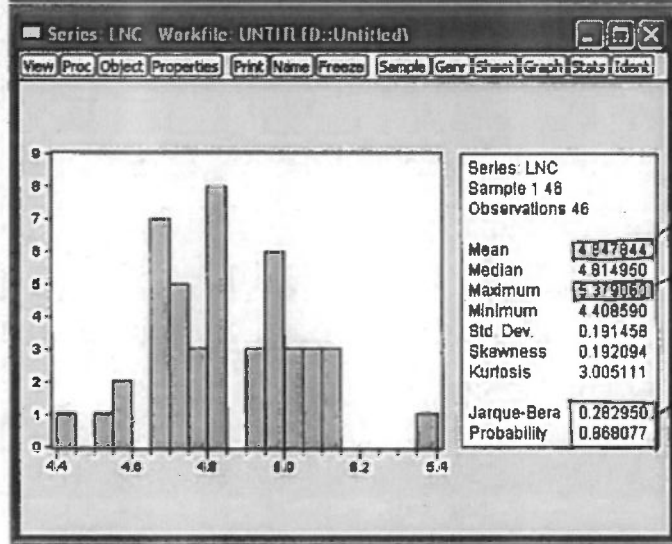
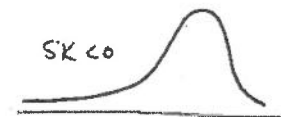
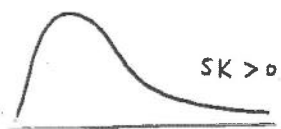
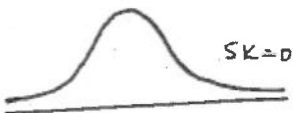
$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

standard dev.

$$SD(x) = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

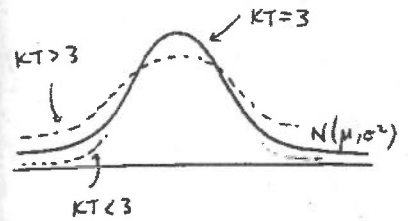
skewness

$$SK = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\right)^{3/2}}$$



kurtosis

$$KT = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4}{\left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\right)^2}$$



127 packs / person

217 packs / person

$$JB = \frac{n}{6} \left[ SK^2 + \frac{(KT-3)^2}{4} \right]$$

$\sim \chi^2(2)$

$H_0: SK=0, KT=3$

$H_1: \neg H_0$

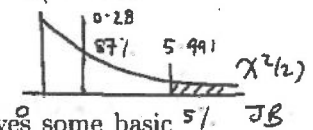


Figure 8: Enter the series name: lnc. The histogram of the data also gives some basic descriptive statistics, up to the standardized fourth central moment (kurtosis), and the Jarque-Bera test for normality of the data: this performs the joint test that the skewness = 0 and the kurtosis = 3.

dnr  $H_0$   
 ("normality")

generally, if "Prob" > 0.05,

do not reject (dnr)  $H_0$  at 95% (etc.)

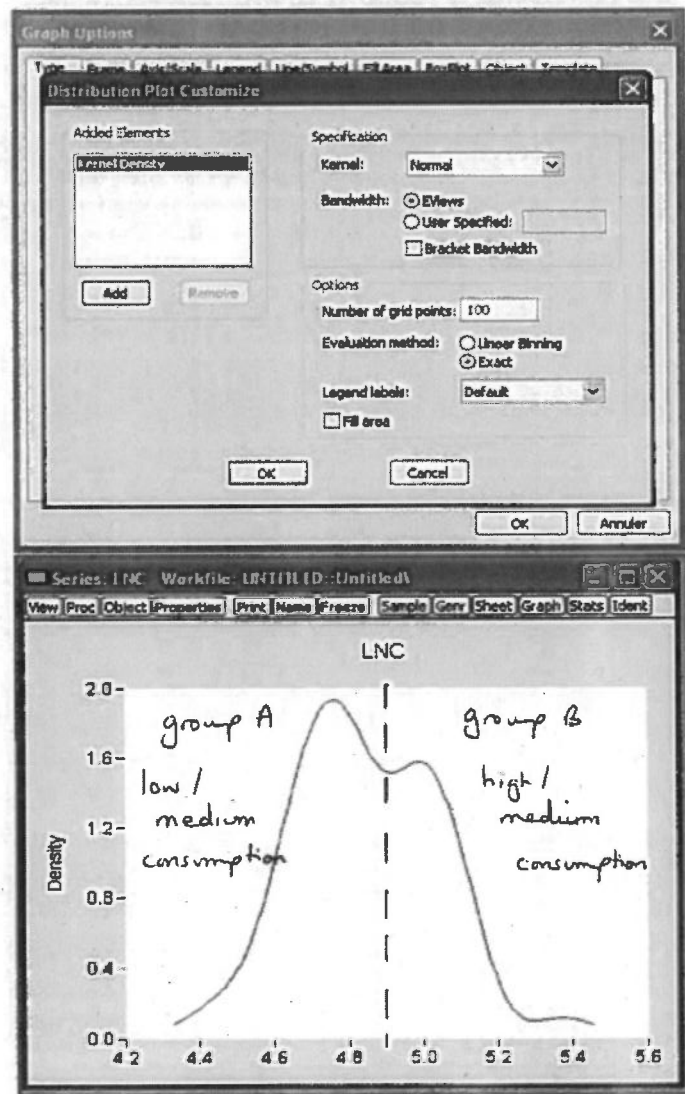
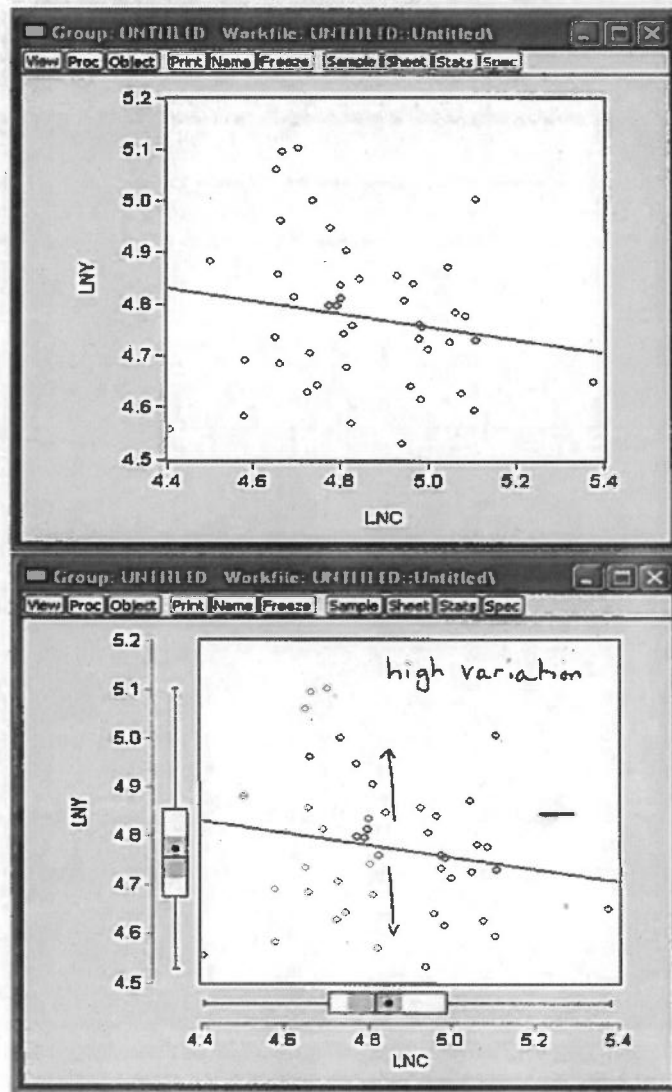


Figure 28: Select 'Exact' instead of 'Linear Binning', to give the Gaussian kernel density plot of the data:  $\hat{f}(x) = (nh)^{-1} \sum_i K((x - X_i)h^{-1})$ , where  $n$  is the sample size,  $h$  is the kernel bandwidth (chosen automatically by EViews), and  $K()$  is the kernel (here, the  $N(0,1)$  pdf). Compare this to the histogram of the data that was plotted earlier.

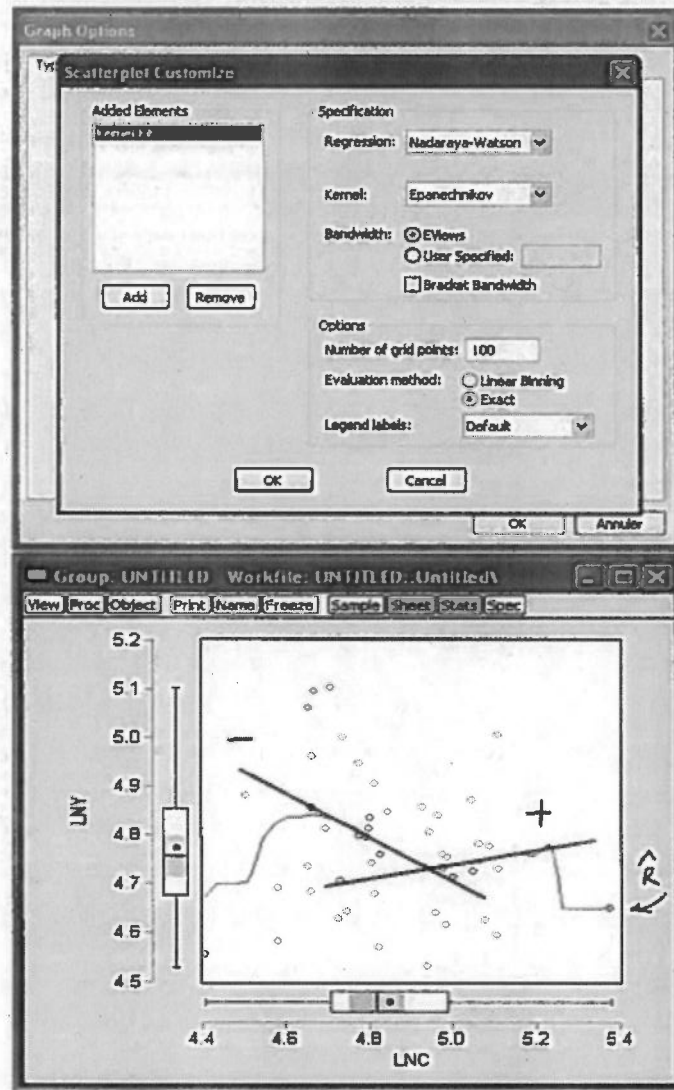


ln Y

suggests that  
 $Y \uparrow \Rightarrow c \downarrow$   
 (strange?!)

Inc

Figure 40: Scatter plot of lny against lnc, with ordinary least squares (OLS) fit from a regression of lny on a constant and lnc. Plot the same figure with boxplots of lny and lnc on the axes!



ln Y

suggests that  
 $Y \uparrow \Rightarrow c \downarrow$  when  
 $c$  low, and  
 $Y \uparrow \Rightarrow c \uparrow$  when  
 $c$  high  
 (strange?!)

ln c

Figure 42: Choose 'Nadaraya-Watson' regression, with 'Epanechnikov' kernel, 'Exact' rather than 'Linear Binning', and bandwidth chosen by 'Eviews', to give the Nadaraya-Watson kernel estimator of  $Y_i$  on  $X_i$  is given by  $\hat{R}(x) = \arg \min_{\psi} \sum_i (Y_i - \psi)^2 K((x - X_i)/h)$ , where  $\psi$  is a locally fit constant, and  $K(u) = (3/4)(1 - u^2)$  on  $[-1, 1]$  is the Epanechnikov kernel.

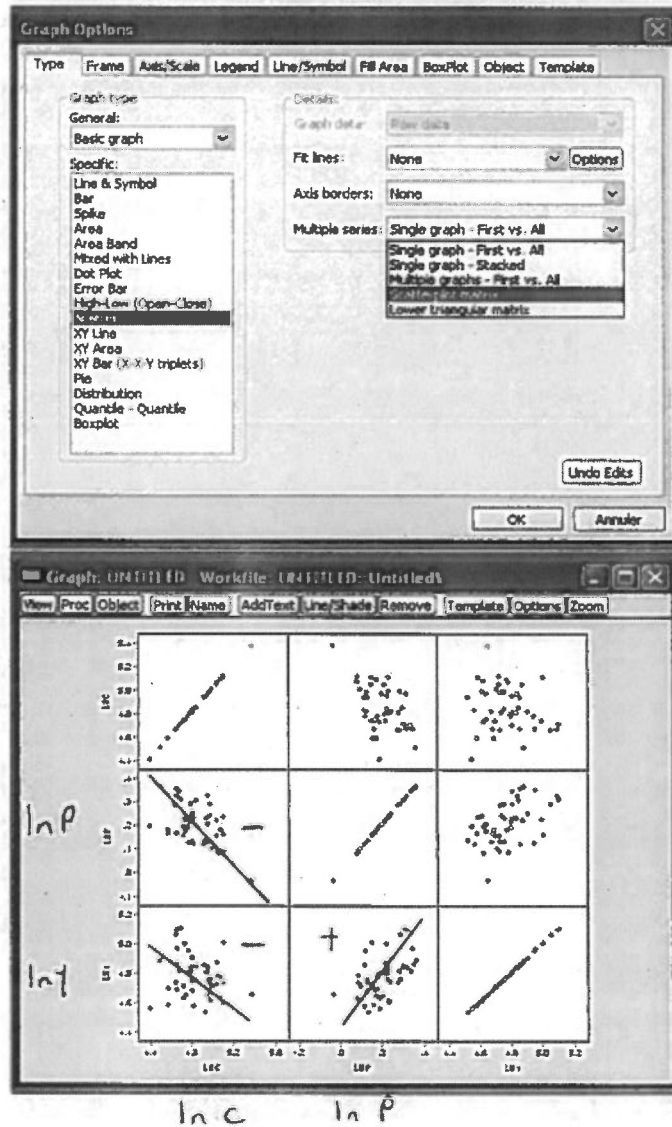


Figure 49: Plot a matrix scatterplot of  $\ln c$ ,  $\ln p$  and  $\ln y$ , and name the graph graph01.

(even if we are only interested in  $Y$ )



$$y = X\beta + u$$

$n \times k$

$$y_i = x_i' \beta + u_i$$

$$i = 1, 2, \dots, n$$

$$n = 46$$

$$k = 2$$

$$y_i = \ln c_i$$

$$x_i = \begin{pmatrix} 1 \\ \ln p_i \end{pmatrix}$$

$$\hat{\beta} = (X'X)^{-1} X'y$$

$k \times 1$

$$\widehat{Var}(\hat{\beta}) = \hat{\sigma}^2 (X'X)^{-1}$$

$$\hat{\sigma}^2 = \frac{\hat{u}'\hat{u}}{n-k}$$

$$\widehat{Var}(\hat{\beta}) = \begin{pmatrix} \widehat{Cov}(\hat{\beta}_1, \hat{\beta}_1) & \widehat{Cov}(\hat{\beta}_1, \hat{\beta}_2) \\ \widehat{Cov}(\hat{\beta}_1, \hat{\beta}_2) & \widehat{Cov}(\hat{\beta}_2, \hat{\beta}_2) \end{pmatrix} = \begin{pmatrix} \widehat{Var}(\hat{\beta}_1) & \widehat{Cov}(\hat{\beta}_1, \hat{\beta}_2) \\ \widehat{Cov}(\hat{\beta}_1, \hat{\beta}_2) & \widehat{Var}(\hat{\beta}_2) \end{pmatrix}$$

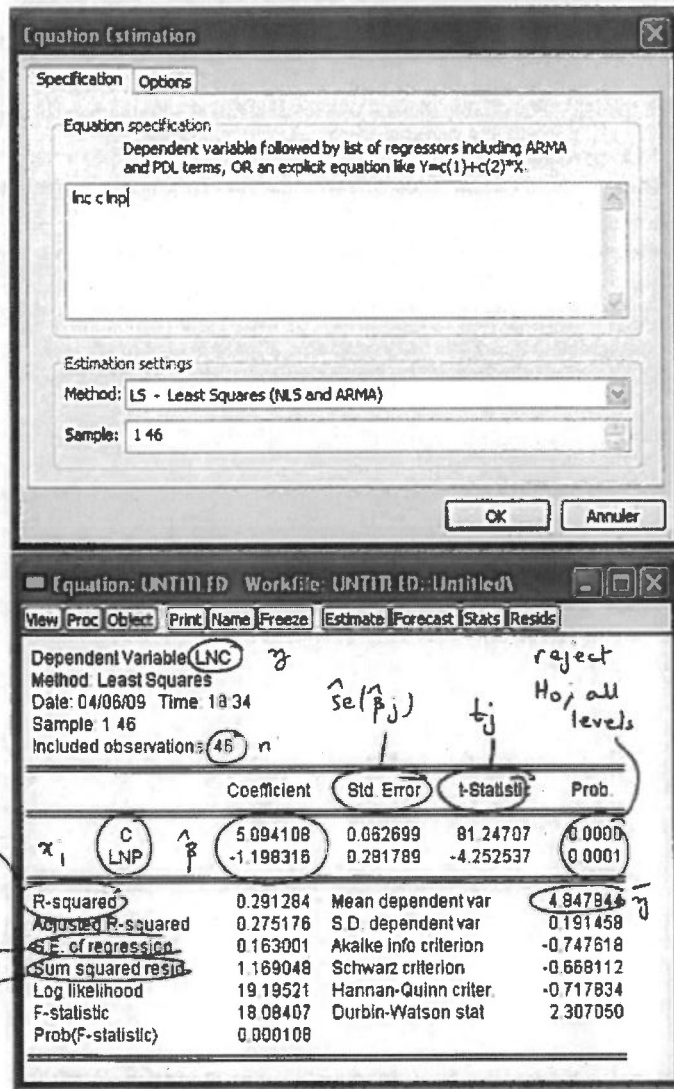
$$\widehat{se}(\hat{\beta}_j) = \sqrt{\widehat{Var}(\hat{\beta}_j)}$$

$$t_j = \frac{\hat{\beta}_j}{\widehat{se}(\hat{\beta}_j)} \sim t(n-k)$$

$$H_0: \beta_j = 0$$

$$H_1: \beta_j \neq 0$$

(test of individual significance)



$R^2$  measures quality of linear fit (0% - 100%)

but need to include constant

in model for this to work

here,  $R^2 \approx 29\%$ . ("good"  $R^2$  depends on context)

Figure 54: (eq01) Enter 'inc c lnp' for a regression of log consumption on a constant and log price. Observe that the regression output gives estimated coefficients (by ordinary least squares), standard errors, probabilities for individual tests of significance, and various statistics, including coefficients of determination, the sum of squared residuals, the Durbin-Watson statistic for autocorrelation, the Akaike and Schwarz Information Criteria, and the F statistic for significance of the entire regression (we will return to these statistics later in the course).

$R^2$  cannot  $\downarrow$  as  $k \uparrow$ , even if irrelevant variables added

adjusted  $R^2$  (denoted  $\bar{R}^2$ ) controls for  $k$ :  $\bar{R}^2 = 1 - \left(\frac{n-1}{n-k}\right)(1-R^2)$

$\Rightarrow \bar{R}^2$  must be used to compare models with different  $k$ .

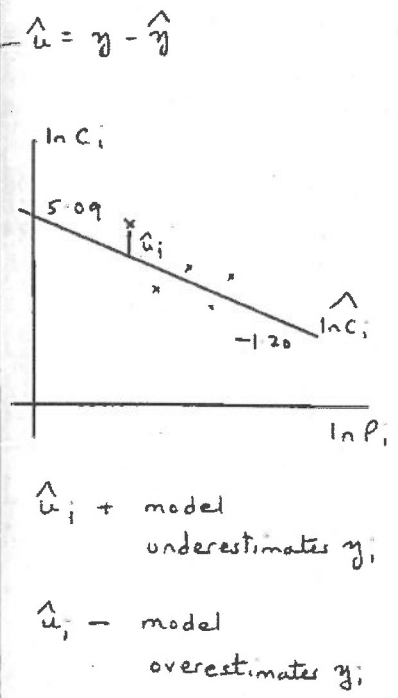
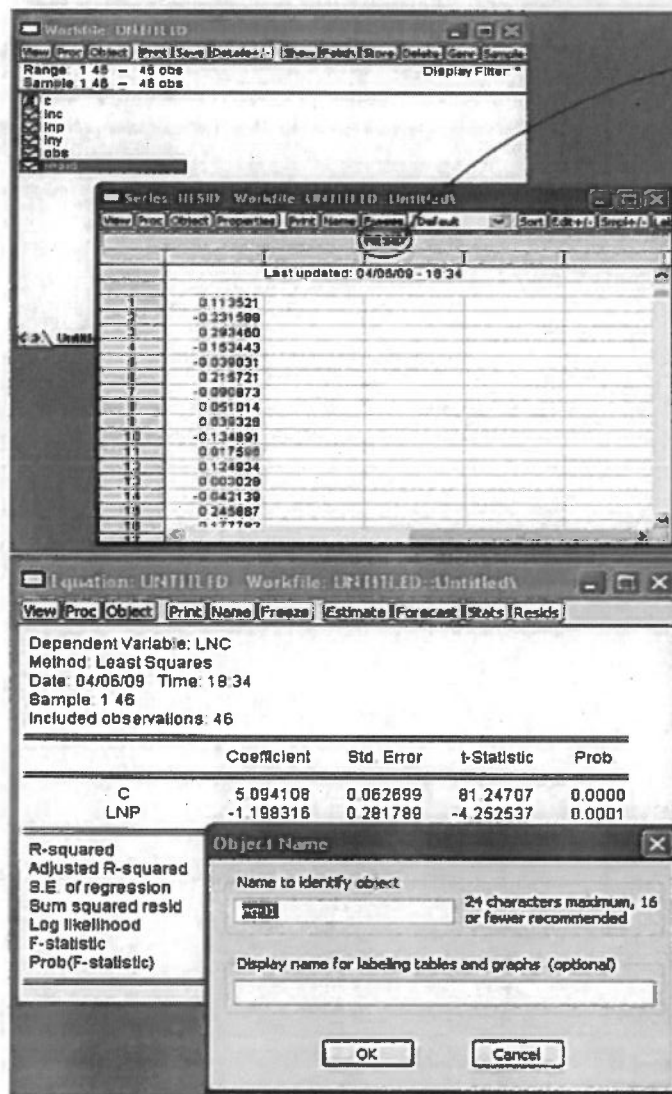


Figure 56: Check that the estimated residuals have been recorded in the workfile object resid: note that the estimated residuals are overwritten each time that a new regression is performed. Name the equation 'eq01'.

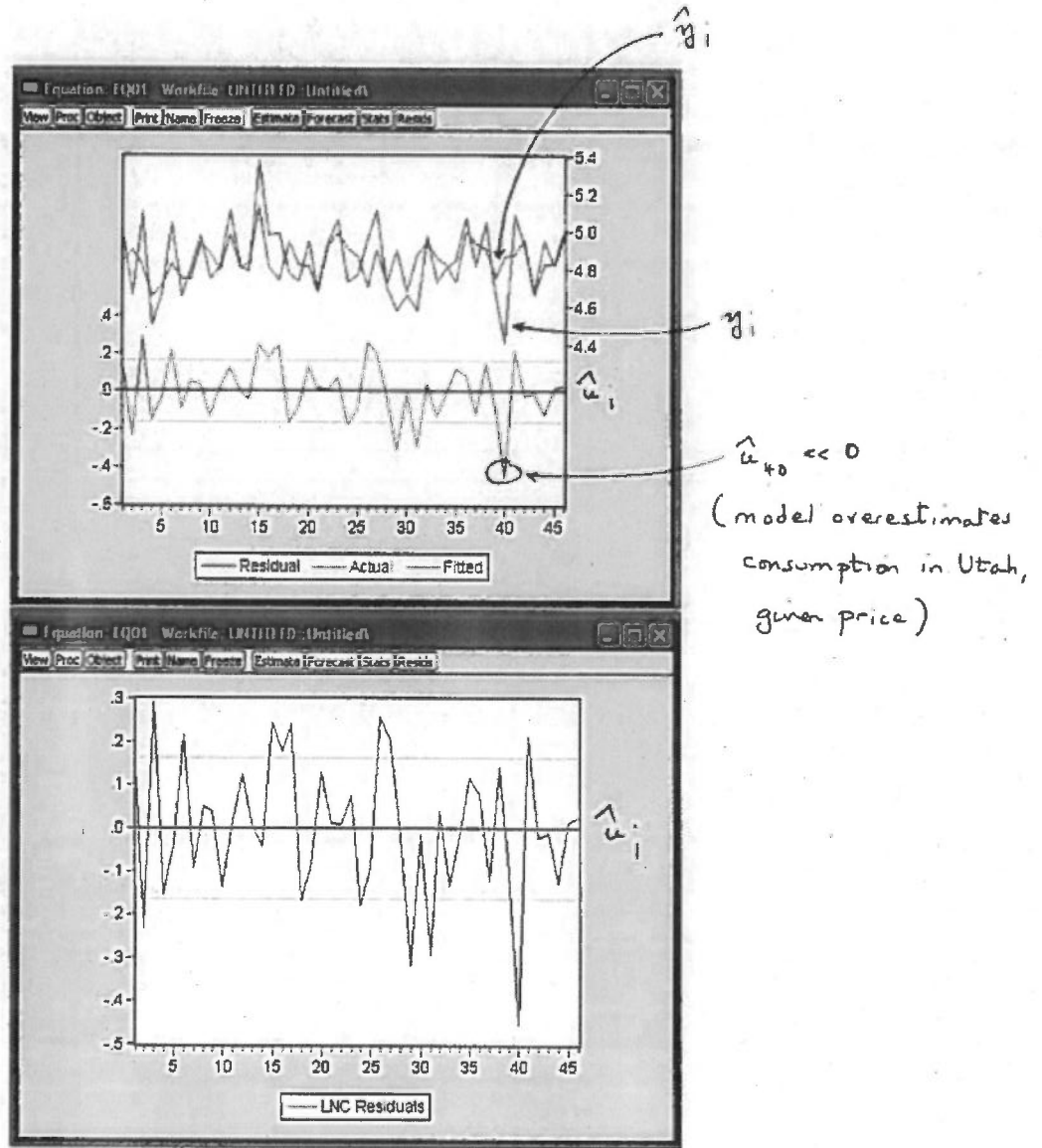


Figure 59: Several representations of the fitted residuals  $\hat{u}$ .

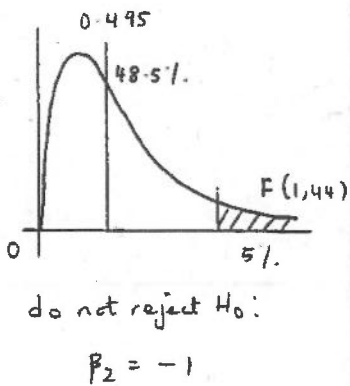
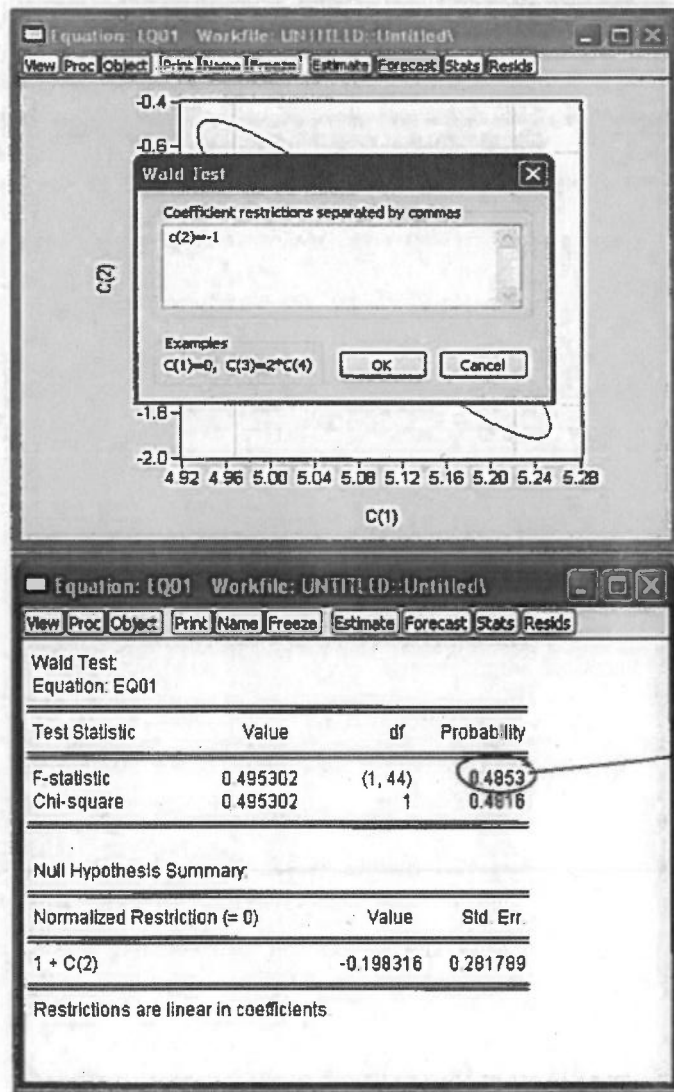
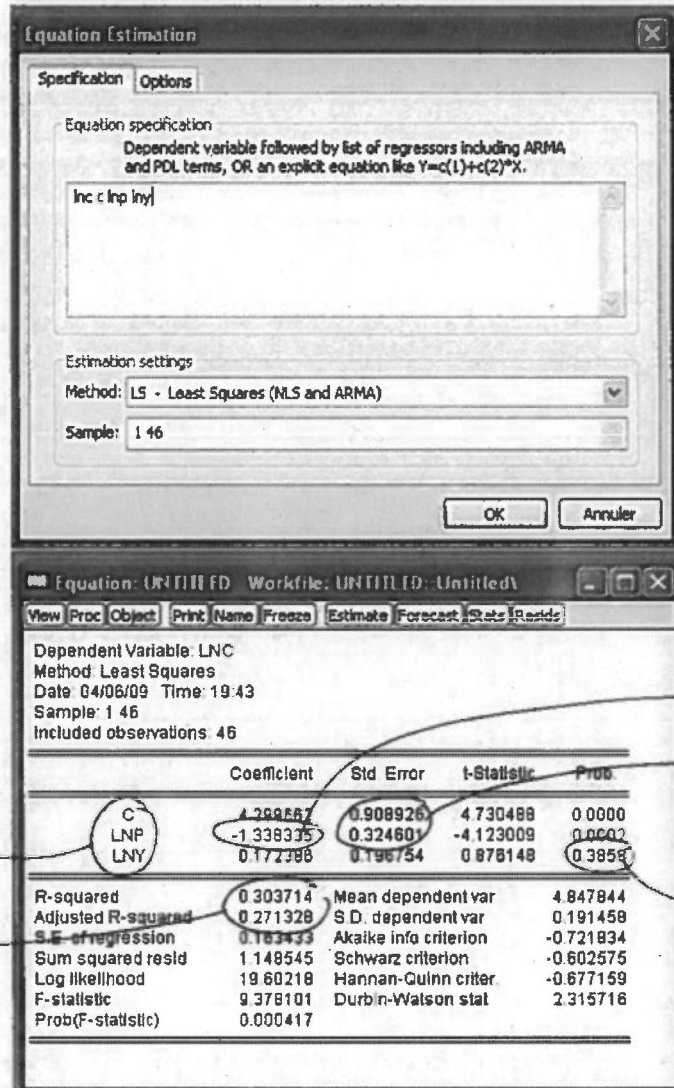


Figure 63: Enter 'c(2)=-1' to perform a Wald test of  $H_0 : \beta_2 = -1$ . Observe that both exact ( $F$ ) and asymptotic ( $\chi^2$ ) statistics are reported: the result of the test is that we do not reject the null at the 90% level (say). Why is an  $F$  test reported, not a  $t$  test?

$$F(1, n-k) = [t(n-k)]^2$$



$C \sim P, Y$

$k=3$

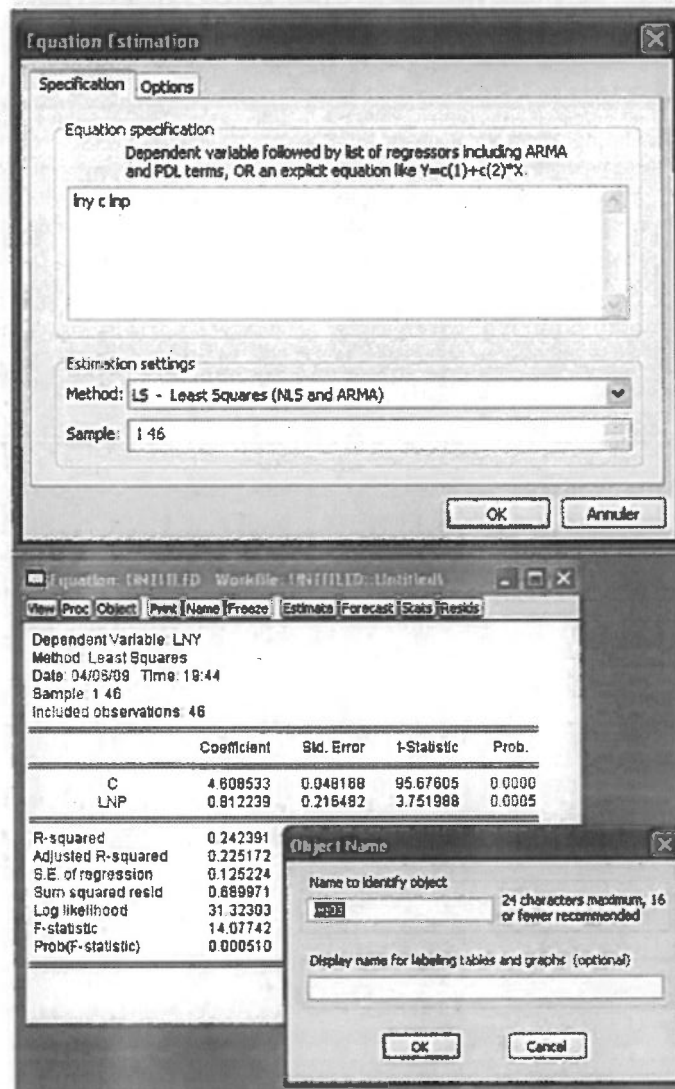
$R^2$  higher than eq01  
 $\bar{R}^2$  lower than eq01

similar to eq01

$\hat{se}$ 's higher than eq01 (have lost precision as  $k \uparrow$ )

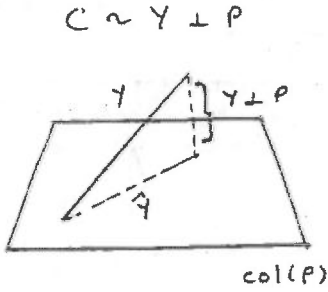
income is not significant

Figure 72: (eq02) Run the regression of log consumption on a constant, log price, and log income, using ordinary least squares, and consider the regression output.



$Y \sim P$   
 correlation  
 $\nrightarrow$  causality

Figure 74: (eq03) Regress log income on a constant and log price, using ordinary least squares, and consider the regression output. Name this equation 'eq03'.



equation: UNTITLED Workfile: UNTITLED - Untitled1

View Proc Object Print Name Freeze Estimate Forecast Stats Resids

Dependent Variable: LNC  
 Method: Least Squares  
 Date: 04/06/09 Time: 19.45  
 Sample: 1 46  
 Included observations: 46

	Coefficient	Std. Error	t-Statistic	Prob.
C	4.847844	0.028370	170.8795	0.0000
RESID_EQ03	0.172386	0.231645	0.744183	0.4607

R-squared: 0.012430  
 Adjusted R-squared: -0.010015  
 S.E. of regression: 0.192414  
 Sum squared resid: 1.629025  
 Log likelihood: 11.56400  
 F-statistic: 0.553808  
 Prob(F-statistic): 0.460722

Object Name

Name to identify object: EQ04 (24 characters maximum, 16 or fewer recommended)

Display name for labeling tables and graphs (optional):

OK Cancel

---

Workfile: UNTITLED

View Proc Object Print Save Details+/- Show Fetch Store Delete Genr Sample

Range: 1 46 -- 46 obs  
 Sample: 1 46 -- 46 obs  
 Display Filter: \*

- c
- eq01
- eq02
- eq03
- eq04
- inc
- inp
- iny
- obs
- resid
- resid\_eq03

Untitled / New Page /

Figure 77: Name this equation 'eq04', and then select the eq01 object.

$C \perp P \sim$   
 $Y \perp P$

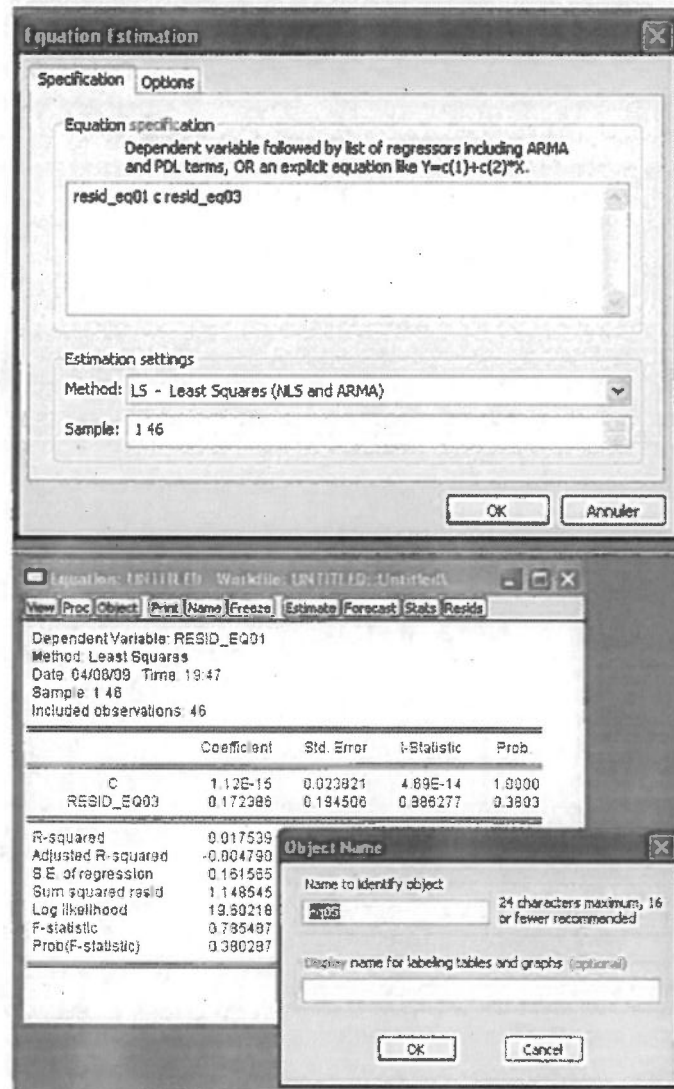
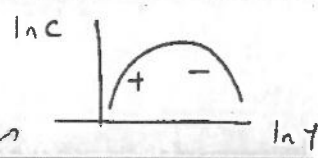


Figure 80: (eq05) Run the regression of resid\_eq01 on a constant and resid\_eq03, and name this equation 'eq05'.

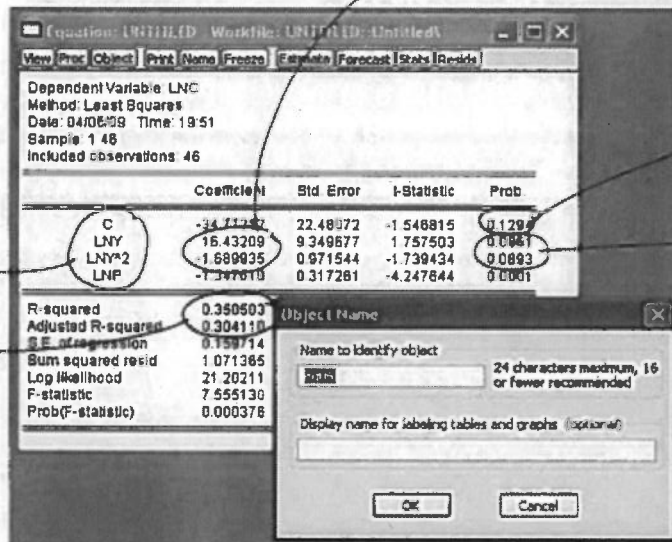




$C \sim Y, Y^2, P$

$k=4$

highest  $\bar{R}^2$  of all models for  $\ln c$ , (eq01, eq02, eq06)



"almost" significant at 10% level

significant at 10% level.

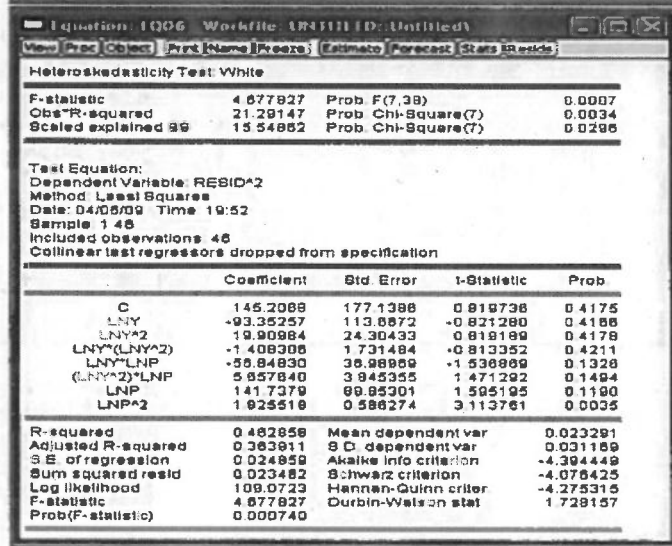


Figure 84: (eq06) Perform the regression of log consumption on a constant, log income, log income squared, and log price, and name the equation 'eq06'. Perform White's  $nR^2$  test for heteroscedasticity: note that homoscedasticity is rejected at all usual levels of significance. (This is evidence against the classical assumption  $Var(u) = \sigma^2 I_n$ ).