# Solution to Task T3.

## Data management in Gretl.

## Task T3.1. Generating Gretl data files

Beach umbrella rental

a. Enter data into Gretl manually.

                    File --> New data set

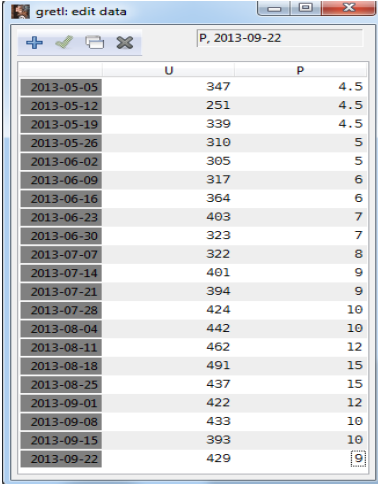Number of observations: 21

Structure of data set: time series

Daily date to represent week: Sunday

Time series frequency: weekly
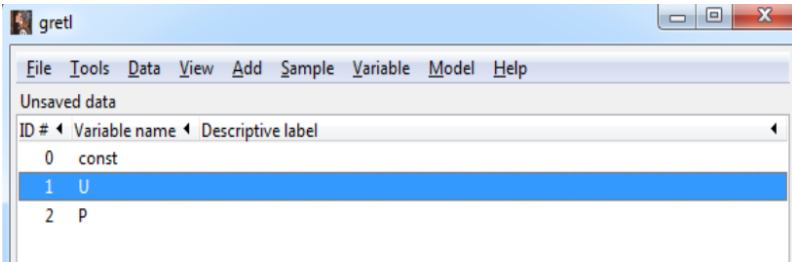
Starting observation: 2013-05-05 (first week of may 2013)

Select: start entering data values.

Give a name to each variable and start entering data values manually.



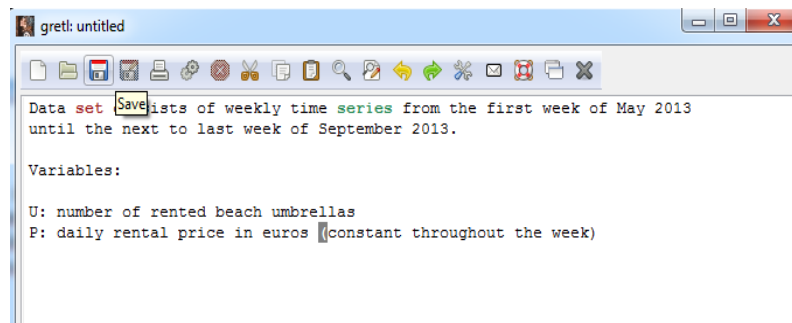The main window shows variables $U$ and $P$.

b. Add a description for each variable.

Highlight the variables in the main window, right-click and select the *Edit attributes* option from the pulldown menu. Write down the description for each variable.



c. Write down some notes on the data set.

Click on the *Session icon view* icon of the Toolbar. Then, click on the *Data info* icon and write down the text shown in the window below.



Before quitting click on the *Save* icon.

d. Save the data file.

```
File --> Save data as ...
```
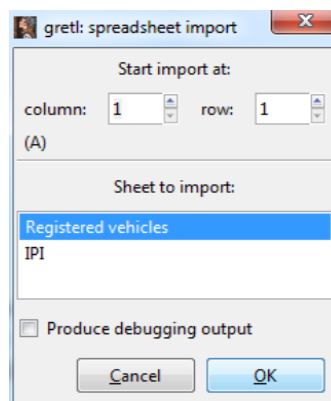
Select the folder where you want to save the data and use the name `umbrellas.gdt`.
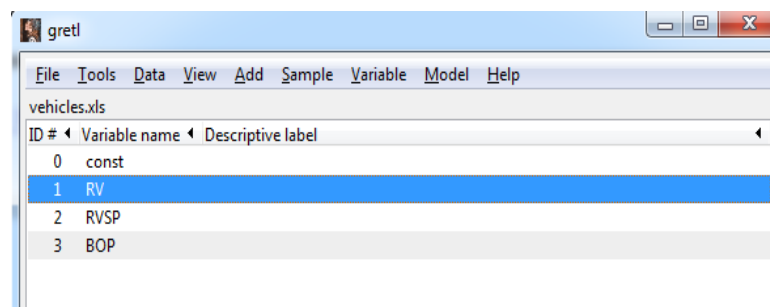
## Registered vehicles

a. Read data from an Excel file.

```
File --> Open data   --> User file
```

Go to the folder where the Excel data file (`vehicles.xls`) is. Load this file into Gretl. In the dialog box, select the sheet, and the row and column where the data start.
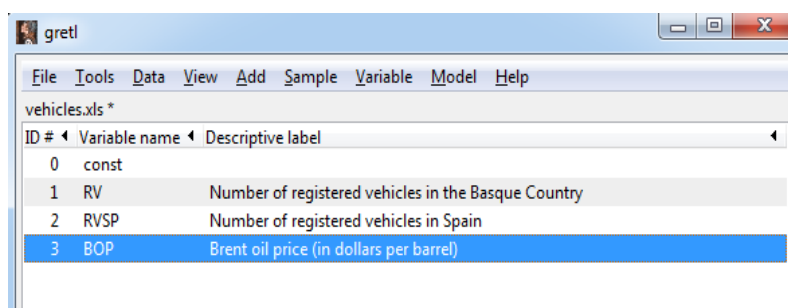
The three imported variables appear in the main window.



b. Add a description for each variable.

Highlight the variables in the main window, right-click and select the *Edit attributes* option. Write down the description for each variable.



c. Save the data file.

```
File --> Save data as ...
```

Go to the folder where you want to save the data and use the name `vehicles.gdt`.

d. Change the data set structure.

```
Data --> Dataset structure
```

Then, mark the <u>Time series</u> option and give the time series frequency and the starting observation.

***DON'T forget to save the changes to this data file***

<u>Wages</u>

a. Save the data file in your folder.

```
File --> Open data   --> Sample file ...
```
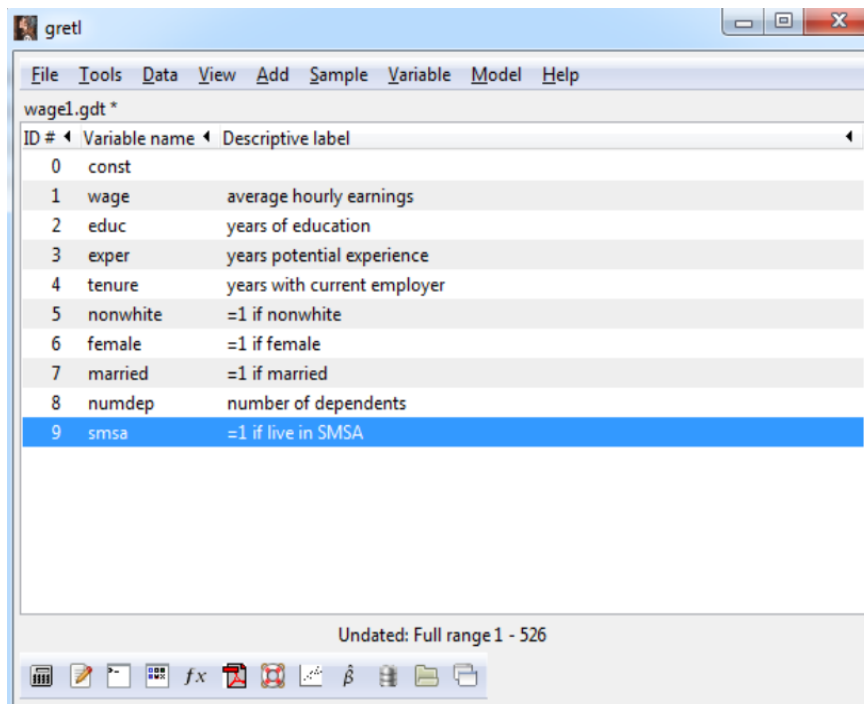
Choose the *Wooldridge* tab and load the file called `wage1`.

Save this data file in your own folder.

```
File --> Save data as ...
```

b. Delete the variables you are not interested in.

Highlight all the variables you are NOT interested in using the cursor. Right-click and select the *Delete* option from the pulldown menu. After having done this, the main window looks as follows:



c. Change the name and the attributes of the variables.

Highlight each variable, right-click and select the option *Edit attributes*. Then, write down the new names and attributes.



d. Do NOT forget to save all the changes in a data file called `wages.gdt`.

```
File --> Save data as ...
```

## Holiday cottages in Biscay

a and b. Load the data file `HCBiscay.txt` and check the codes that Gretl assigns to the qualitative variables.
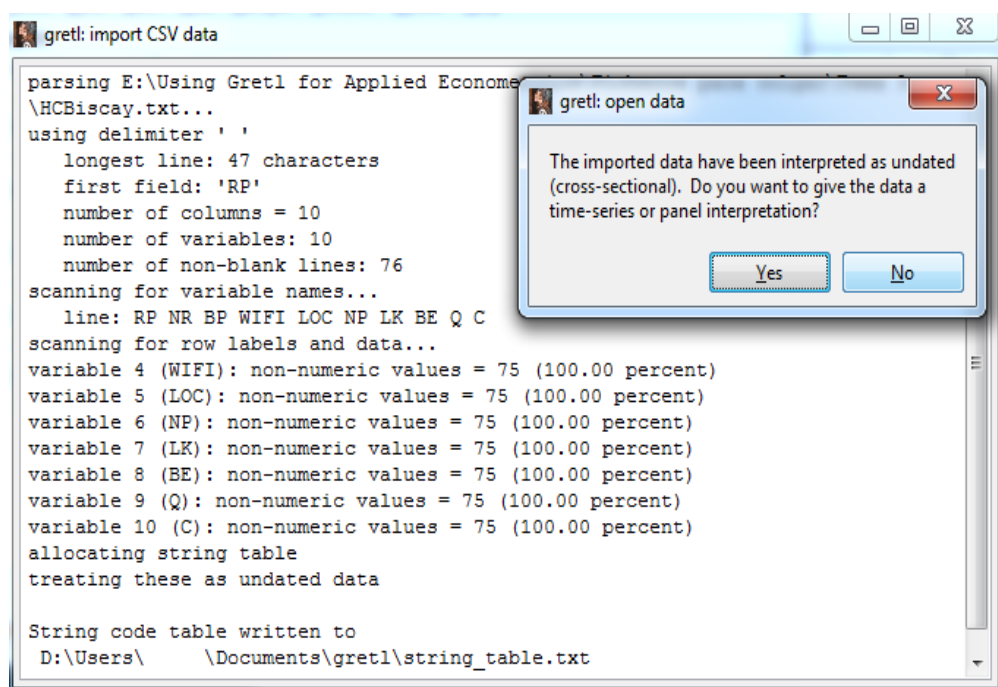
<div align="center">

`File --> Open data    --> User file`

</div>

Look for the folder where the data file `HCBiscay.txt` is and load the data file.

When a data file in text format is loaded, Gretl generates two windows that provide all the data information.
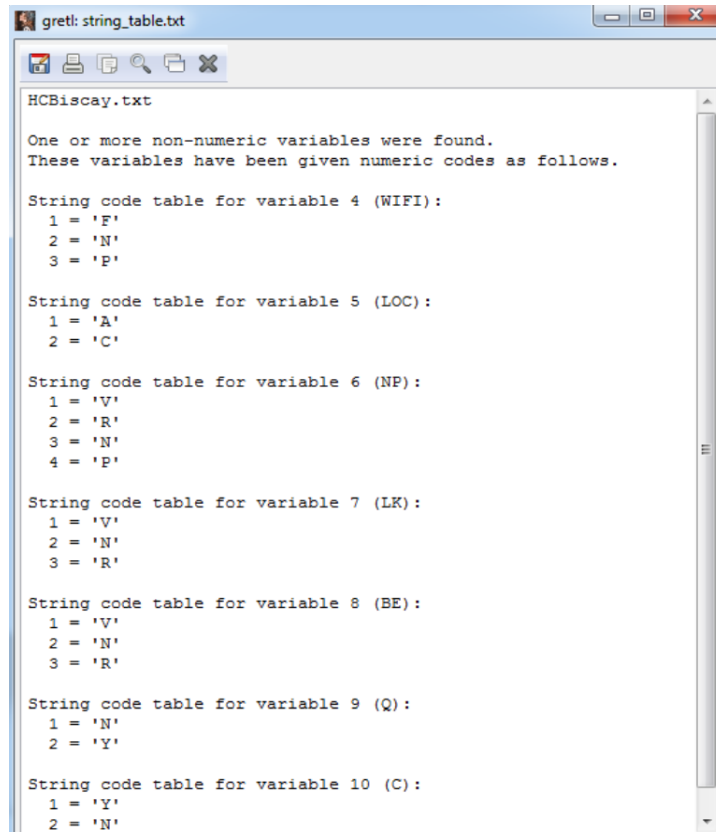
The window **gretl: import CSV data** shows:

- The name of the data file imported: HCBiscay.txt.

- The number of observations: 76

- The qualitative variables with non-numeric values: WIFI, LOC, NP, LK, BE, Q and C.

- At the bottom, you can find the name of the file where the codes for the qualitative variables, that is, the numeric values assigned by Gretl to the categories of the qualitative variables are: string_table.txt.
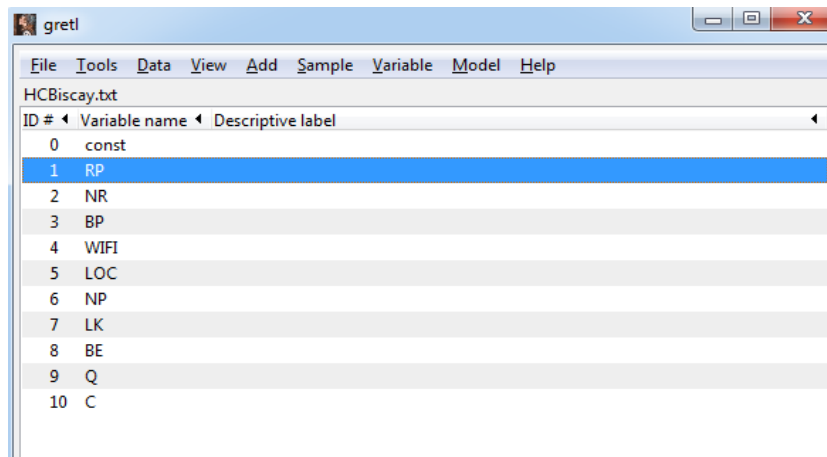


The window **gretl: string_table.txt** shows the codes assigned to the categories of the qualitative variables. For instance, the label $A$ assigned in the text file `HCBiscay.txt` to a holiday cottage far from the town center is transformed into the number 1 in the Gretl format file.

Note that the same numerical value is not always assigned to the same label. For instance, the label $R$ becomes the numerical value 2 for the qualitative variable $NP$, and the numerical value 3 for $LK$. This is due to the fact that Gretl always assigns the numerical value 1 to the first label that detects for each variable.

All the imported variables, quantitative and qualitative, are listed in the main window of Gretl.
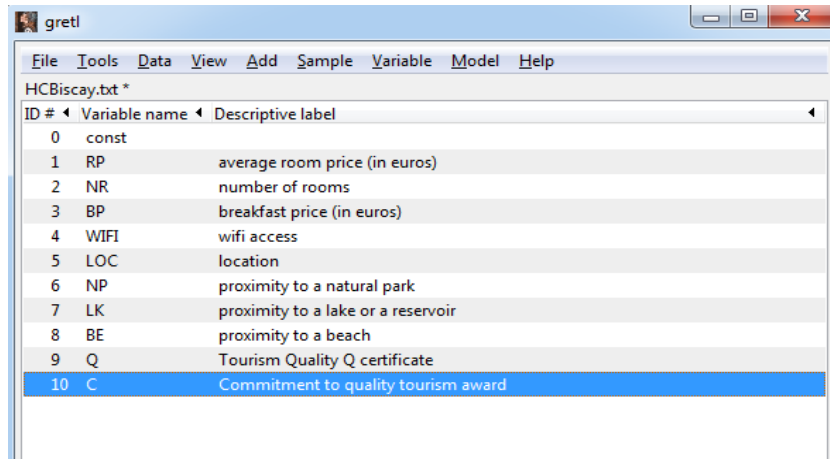


Highlight ONE of the variables, right-click and select the option *Display values* from the pulldown menu to see the codes used in the initial file `HCBiscay.txt`.
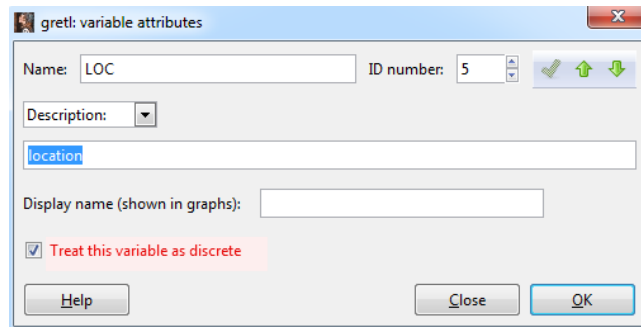
Highlight ALL the variables, right-click and select the option *Display values* to see the numerical values assigned by Gretl.

c. Add a description for each variable and verify that all the qualitative variables have been defined as discrete variables.

Highlight each variable in the main window, right-click and select the *Edit attributes* option. Write down the description of the variable.

Let's verify whether the qualitative variable *LOC* has been defined as a discrete variable. Highlight the variable *LOC* using the cursor, right-click and select the option *Edit attributes*. As you may see in the window below, the option <u>Treat this variable as discrete</u> is on.



Repeat this commands for the rest of the qualitative variables.

d. Do NOT forget to save all the changes in a data file called `cottages.gdt`.
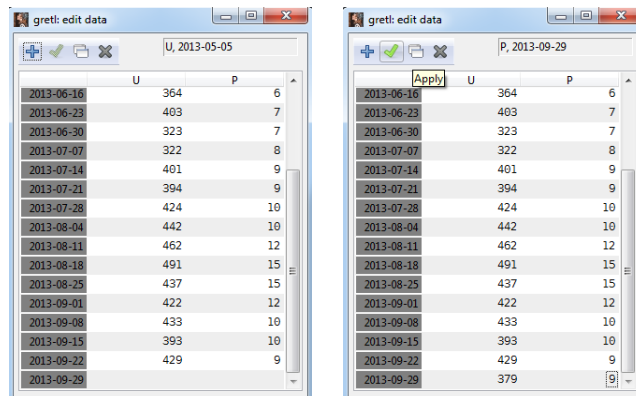
# Task T3.2. Modifying Gretl data files

Beach umbrella rental

I. Add new observations manually. Click

```
Data --> Add observations ...
```

and write down the number of observations to add. Then, click

```
Data --> Edit values
```

and write down the new values in the empty row.



II. Add variables manually.

```
Data --> Edit values
```

Click on the symbol $+$ , select *Add variable* and write down the sample values for the two variables.



Follow the procedure explained in Task T3.1 to edit the attributes for the two new varaibles.

*Do NOT forget to save all the changes into your data file.*

<div style="text-align: center; border: 1px solid blue; display: inline-block;">Registered vehicles</div>

a. Import new data.

                    File --> Append data    --> Excel

Open the folder where the Excel file is and load it. In the dialog box, select the sheet where the new data are and the row and column where the data start.
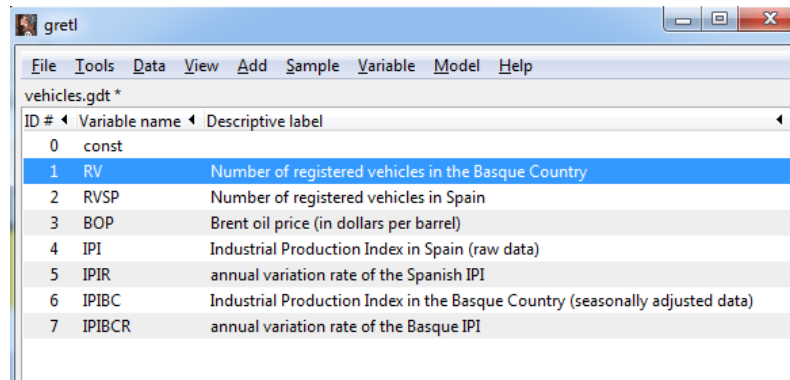


The four imported variables are listed in the main window.

b. Add a description for each new variable.

Highlight each variable using the cursor, right-click and select the option *Edit attributes*. Write down the description of each variable.



c. Do NOT forget to save all the changes into the data file `vehicles.gdt`.

### Simulation

To create a new data set, click

```
File --> New data set
```

Number of observations: 1000
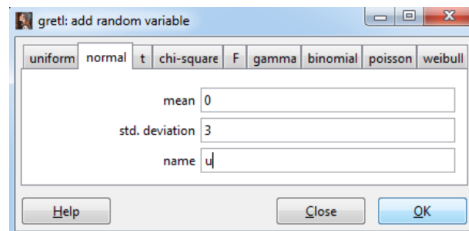
Structure of data set: cross-sectional

DO NOT mark start entering data values, because the data are going to be simulated.

Only the variable index is shown in the main window.



a y b. Simulate the variables:

```
Add -- > Random variable ...
```
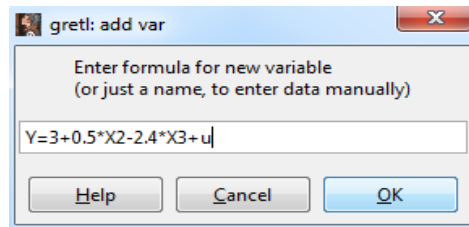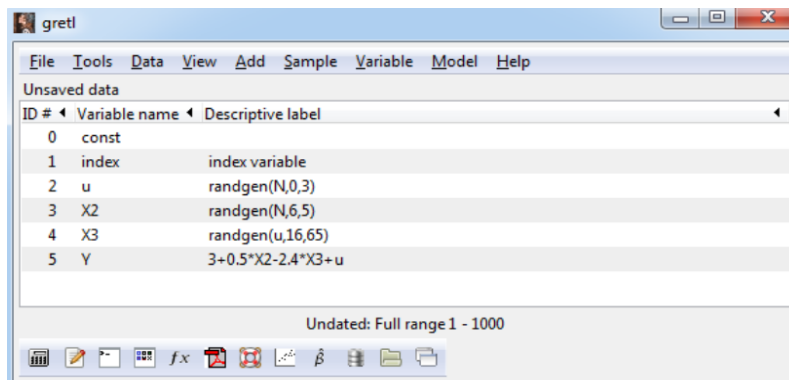
N(0, 9)

N(6, 25)

U(16, 65)

$$Y = 3 + 0,5X_2 - 2,4X_3 + u$$

Now, all the generated variables are listed in the main window.



d. Do NOT forget to save all the changes in a data file called `simulation.gdt`.
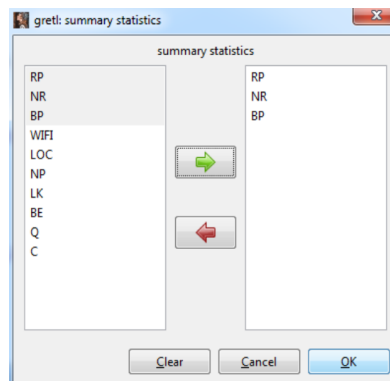
# Task T3.3. Data analysis using Gretl

## Holiday cottages in Biscay

a. Descriptive statistics.

```
View --> Summary statistics
```

Select the quantitative variables: price of a room, number of rooms and price of breakfast.
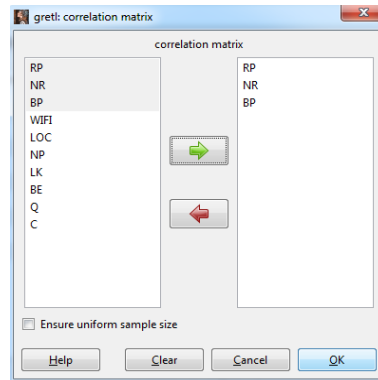


The results are shown in the window below.



- Range of the average price of a room: 133.75 - 25.92 = 107.83.

- Range of the number of rooms: 6 - 1 = 5.

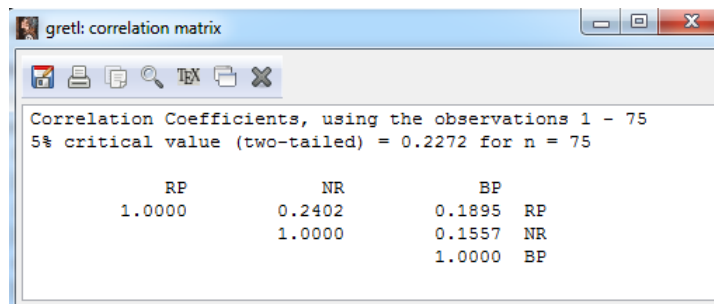- Minimum price of breakfast: €0.

- Maximum price of breakfast: €9.

b. Correlation matrix.

```
View --> Correlation matrix
```

Select the quantitative variables: price of a room, number of rooms and price of breakfast.

The output window shows that the pairwise correlation between the variables is positive but quite small.
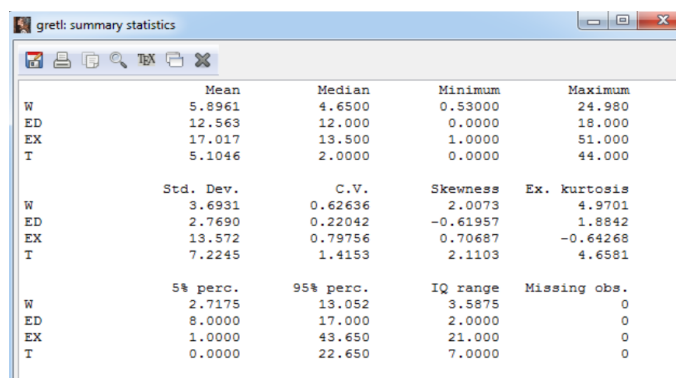


## Wages

a. Descriptive statistics.

                    View --> Summary statistics

Select the variables wage, education, experience and tenure. The results appear in the window below.
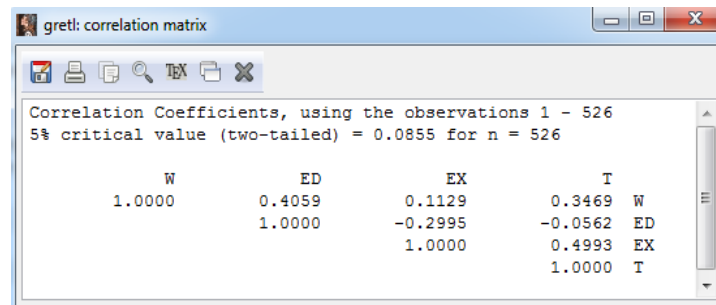


This table shows the main descriptive statistics for the selected variables. Note that the variables that show more variability are experience and tenure.

b. Correlation matrix.

<div align="center">

`View --> Correlation matrix`

</div>

Select the variables wage, education, experience and tenure. The results are shown in the window below.



The variable wage is positively correlated with the variables education, experience and tenure. The correlation coefficient between wage and experience is quite small.

Furthermore, education is negatively correlated with experience, the correlation coefficient between education and tenure is close to zero, and the correlation between experience and tenure is strong and positive.
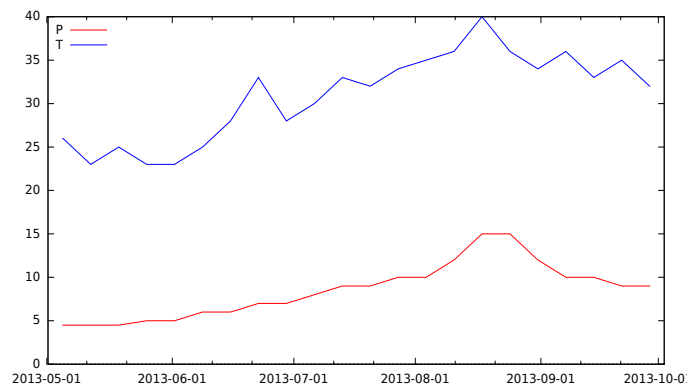
# Task T3.4. Graphical analysis using Gretl.

# Time series graphs.

Beach umbrella rental

To draw several time series on a single graph, click

```
    View --> Graph specified vars      --> Time series plot ...
```

Select the variables price and temperature.





It may be observed that:

- Both series show an increasing trend until the middle of August when they start decreasing until the end of the sample.
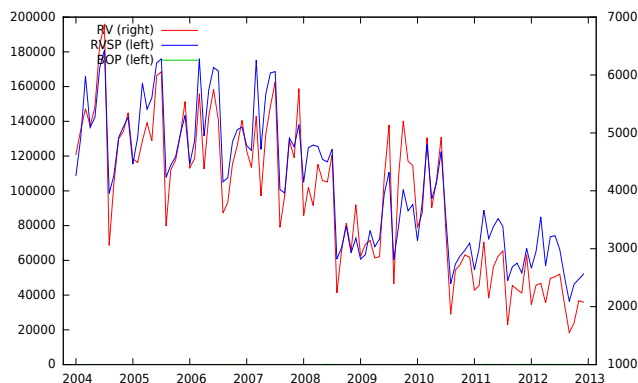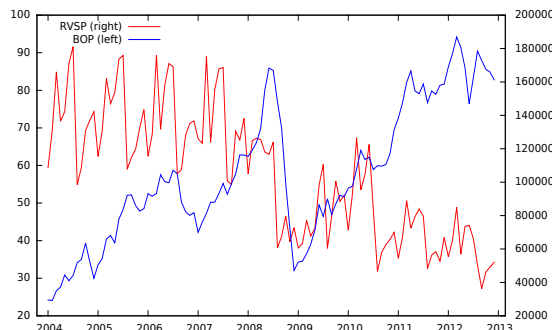
- The series do not present seasonal behaviour.

## Registered vehicles

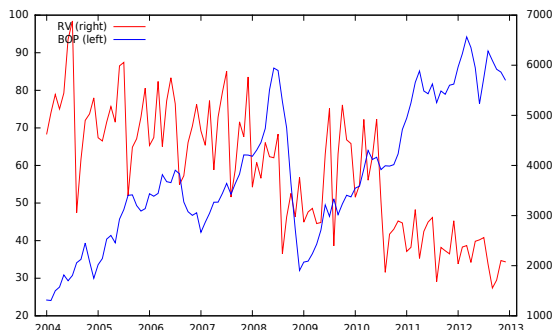a. Vehicles registered and price of Brent.

To draw several time series on a single graph, click

        View --> Graph specified vars        --> Time series plot ...
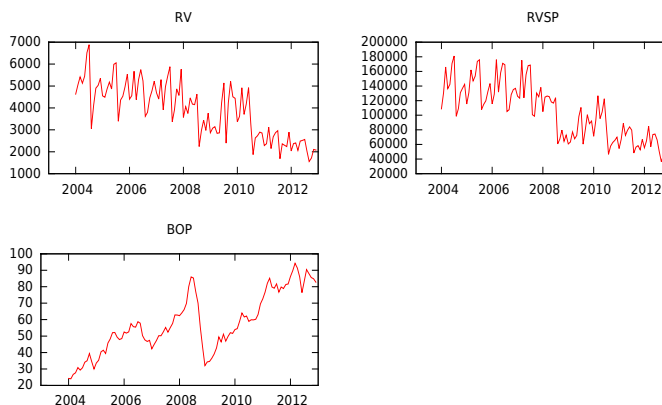
Select the variables of interest.



It is not possible to see the series $BOP$ in this graph because the scale of $BOP$ is much smaller than the scale of the number of vehicles registered in Spain and in the Basque Country. Therefore, in this case, it is better to draw two plots: one of $RV$ and $BOP$ and another one of $RVSP$ and $BOP$



To draw several time series in separate small graphs, click

                View --> Multiple graphs        --> Time series ...
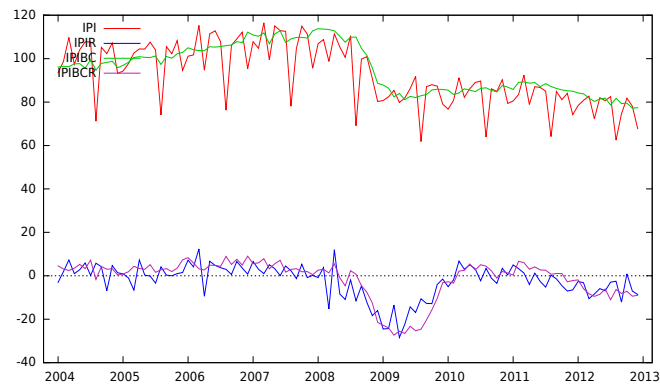
and select the variables of interest.

b. Industrial Production Index.

To draw the series related to IPI on a single graph, click

```
    View --> Graph specified vars    --> Time series plot ...
```

and select the variables of interest.
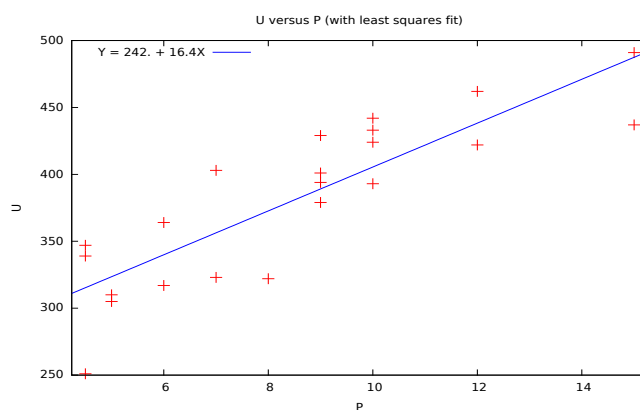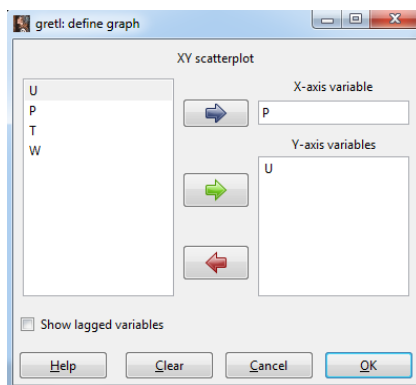


c. Results.

- Trend. The number of registered vehicles series show both a decreasing trend while the trend of the price of Brent series is positive. The IPI series show a slightly increasing trend up to 2008 when they start decreasing reaching negative values. The recession suffered from 2008 to 2010 may be clearly seen in the graph. The annual variation rates turn to be positive at the beginning of 2010, although they became negative again in 2012.

- Seasonality. The series of registered vehicles (both in Spain and in the Basque Country) and the Spanish IPI show seasonal behaviour. It was not expected to observe a seasonal component neither in the IPI of the Basque Country because they are seasonally adjusted data nor in the variation rates because they are annual variation rates.
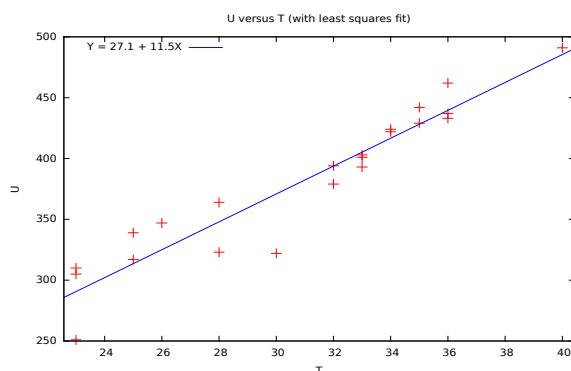
# Scatterplots.

Beach umbrella rental

```
View --> Graph specified vars    --> X-Y scatter ...
```

Select the variable in the Y-axis (umbrella, $U$) and in the X-axis (price, $P$).
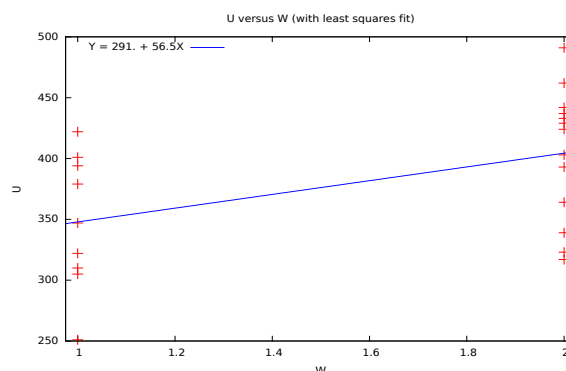




The graph shows a positive relationship between the number of rented umbrellas and the price.

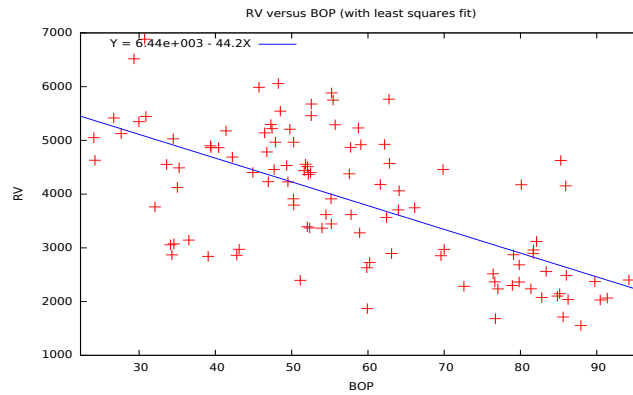Umbrellas against temperature          Umbrellas against $W$





It may be observed in these graphs that there is a positive relationship between the number of rented umbrellas and the temperature, and the number of umbrellas rented and $W$. That is, the higher the temperature is (the less wind there is), the more umbrellas are rented.

## Registered vehicles

```
View --> Graph specified vars    --> X-Y scatter ...
```
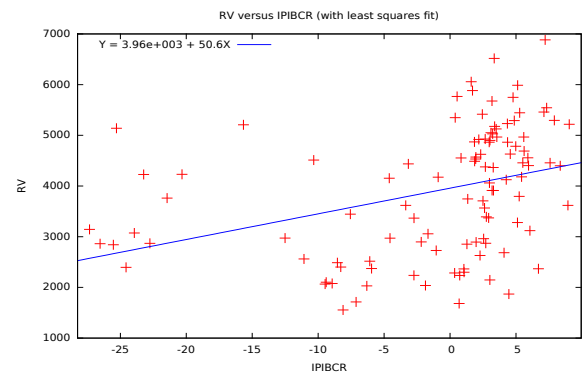
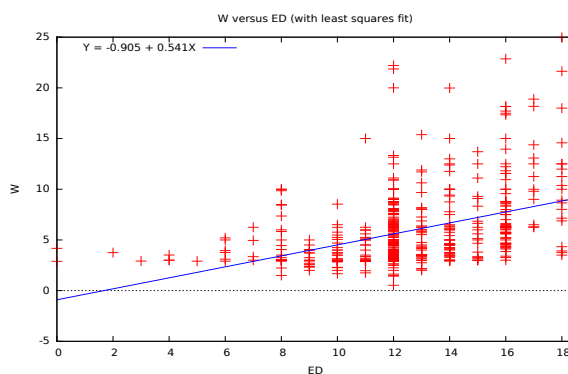Select the variables in the Y-axis (vehicles, $RV$) and in the X-axis (price of oil, $BOP$).



RV versus BOP (with least squares fit)
Y = 6.44e+003 - 44.2X

$RV$ against IPIBC                                      $RV$ against IPIBCR



RV versus IPIBC (with least squares fit)
Y = -2.57e+003 + 67.8X



RV versus IPIBCR (with least squares fit)
Y = 3.96e+003 + 50.6X

These graphs show that there is negative relationship between the number of vehicles registered and the price of Brent, while the relationship between the number of vehicles registered and the IPI (both in raw data and in annual rates) is positive. The IPI indicator is usually used as a proxy to measure the level of economic activity. The higher the economic activity is, the more vehicles are registered.
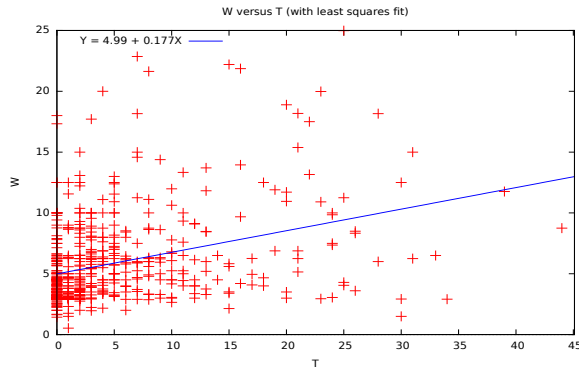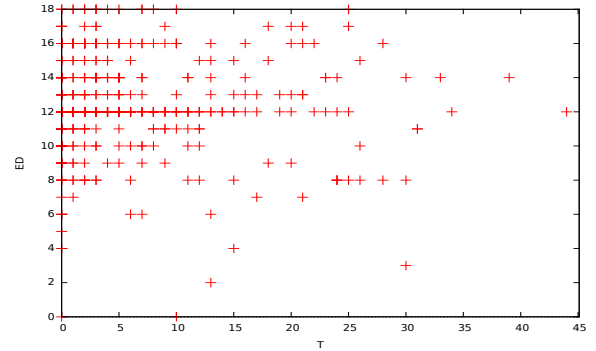
## Wages

*wages* against *education*                              *wages* against *experience*
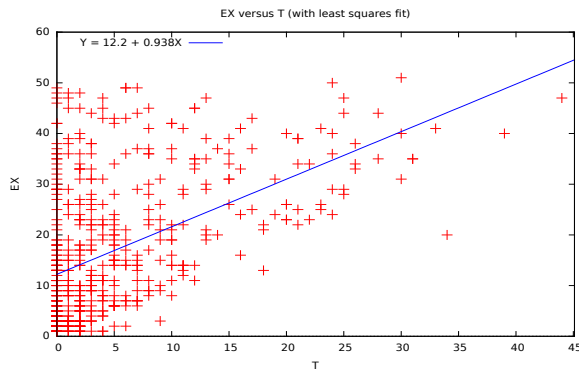


W versus ED (with least squares fit)
Y = -0.905 + 0.541X



W versus EX (with least squares fit)
Y = 5.37 + 0.0307X

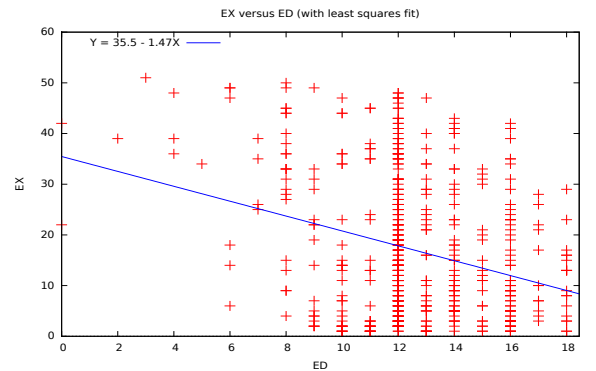### *wages* against *tenure*



### *education* against *tenure*



### *experience* against *tenure*



### *experience* against *education*



Analysing the graphs, it may be concluded that there is a positive relationship between wages and education and between wages and tenure. The larger the education/tenure is, the higher the wages are. However, the relationship between wages and experience does not seem to be very strong.

There is a strong and positive relationship between experience and tenure, whereas the relationship between education and experience or tenure is quite weak.

# 3D graphs.

Beach umbrella rental

To obtain a 3D graph, click:

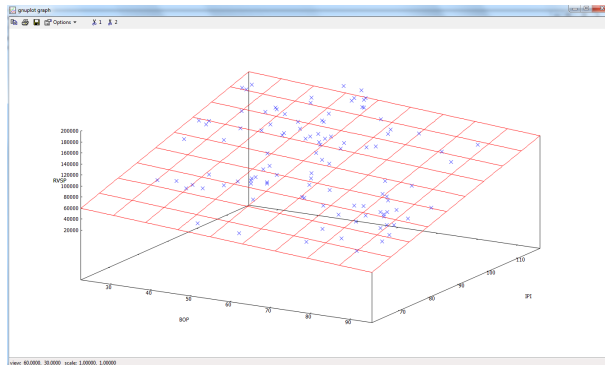```
View --> Graph specified vars    --> 3D plot...
```

Select the variables for the Z-axis (number of umbrellas, $U$), for the X-axis (price, $P$) and for the Y-axis (temperature, $T$).
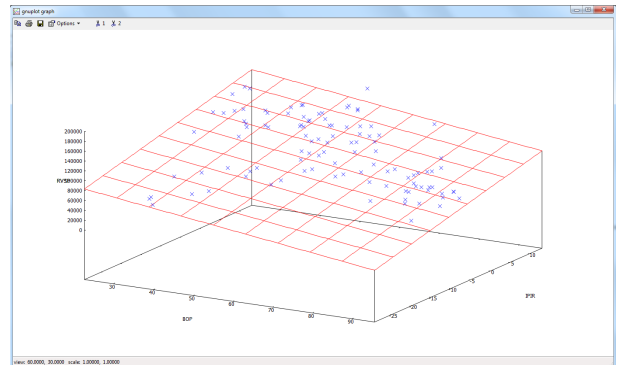




It may be observed that the relationship between number of umbrellas and temperature is strong and positive and the relationship between number of umbrellas and price is not so strong and negative.
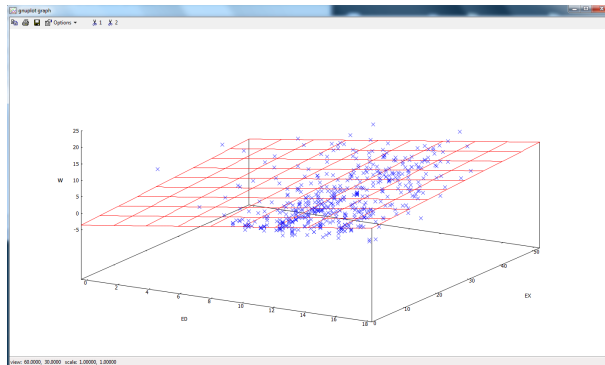
## Registered vehicles

*RVSP* against *BOP* and *IPI*



*RVSP* against *BOP* and *IPIR*



It may observed in the graph that if the price of Brent increases the number of registered vehicles falls, while if the IPI (or its annual variation rate) increases, the number of registered vehicles goes up.
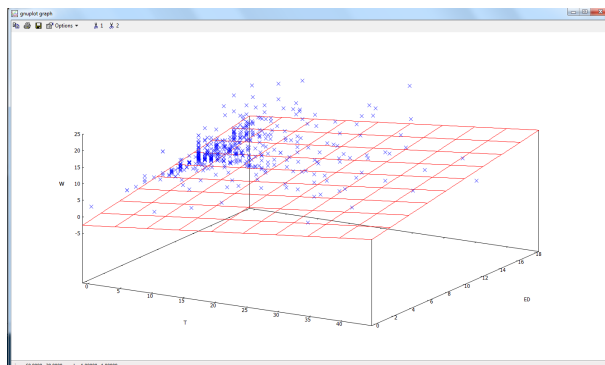
## Wages
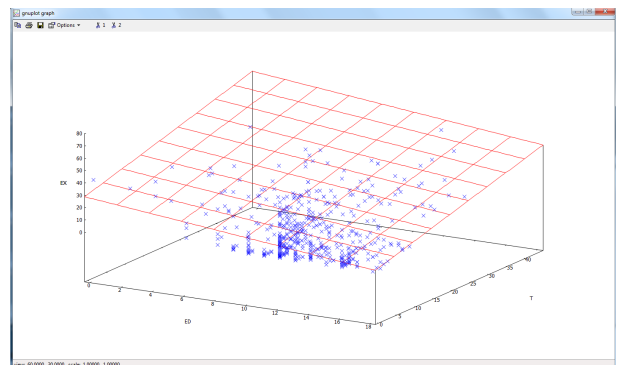
*wages* against *education* and *experience*



*wages* against *experience* and *tenure*



*wages* against *tenure* and *education*



*experience* against *education* and *tenure*



It may be concluded that the relationship between wages and education or tenure is positive, while the relationship between wages and experience is quite weak.
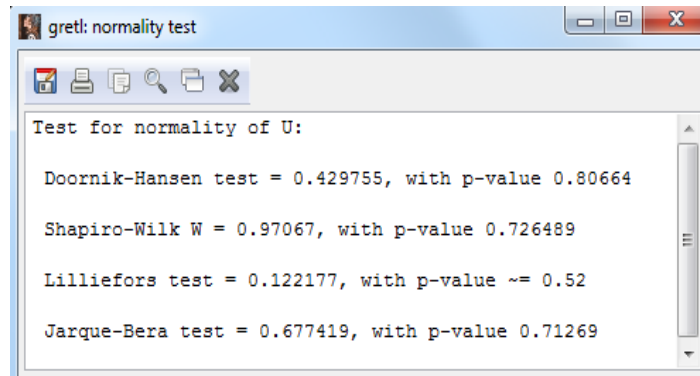
# Task T3.5. Univariate analysis.

Beach umbrella rentals

a. Normality tests.

Select the variable and click

```
Variable --> Normality test
```

Rented umbrellas


Price


Temperature


b. Frequency distribution.
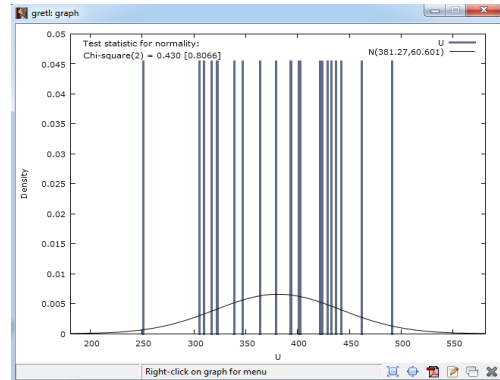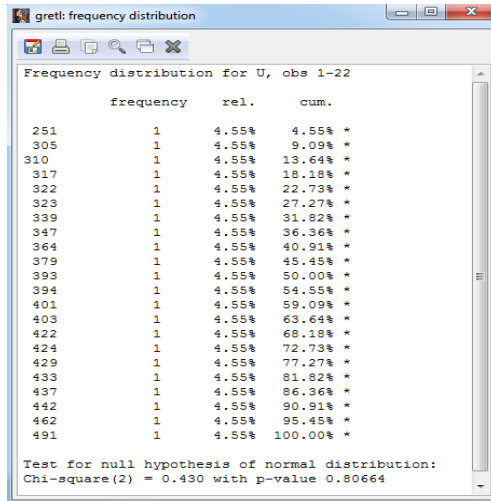
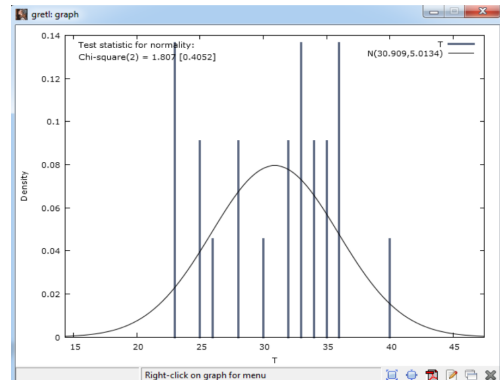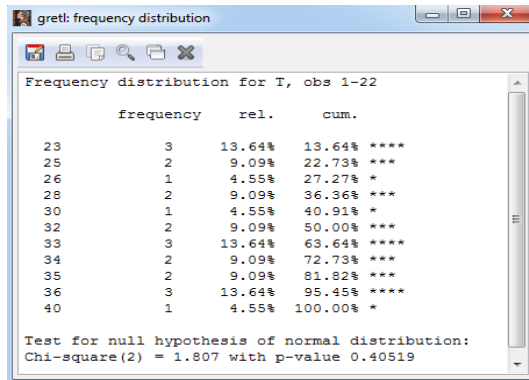To obtain the frequency distribution and the result of the normality test, click

```
Variable --> Frequency distribution ...
```

and mark Test against normal distribution and show plot.
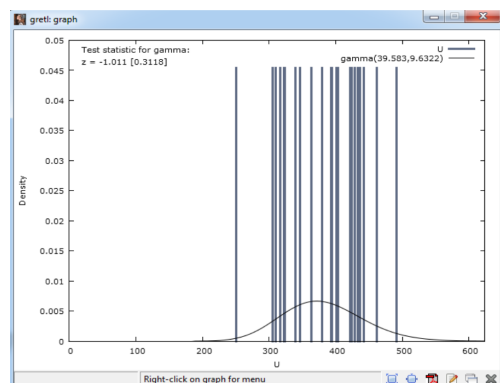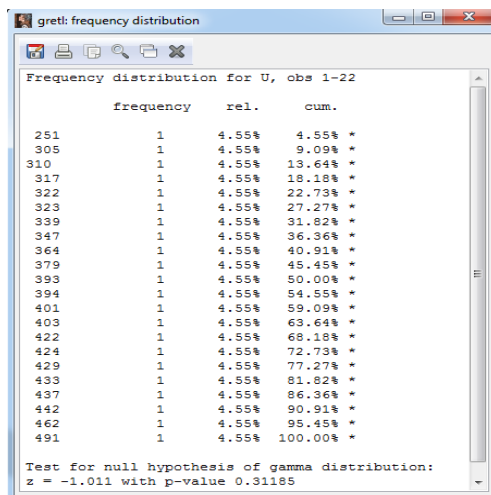
Rented umbrellas



Temperature



To obtain the frequency distribution and the result of the gamma distribution test, highlight the variable of interest using the cursor and click
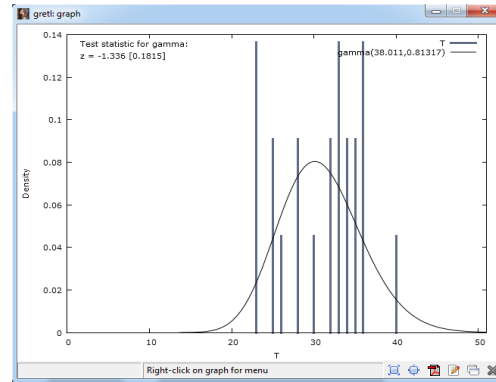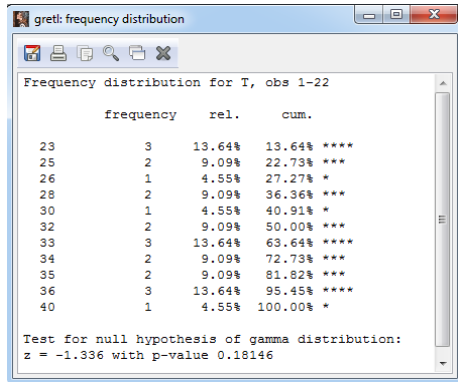
```
Variable --> Frequency distribution ...
```

and mark Test against gamma distribution and show plot.

Rented umbrellas

Temperature



c. Results.

- At a 5 % significance level, the null hypothesis of a normal distribution is not rejected neither for the variable $U$ nor for the variable $T$.

- At a 5 % significance level, the null hypothesis of a gamma distribution is not rejected neither for the variable $U$ nor for the variable $T$.
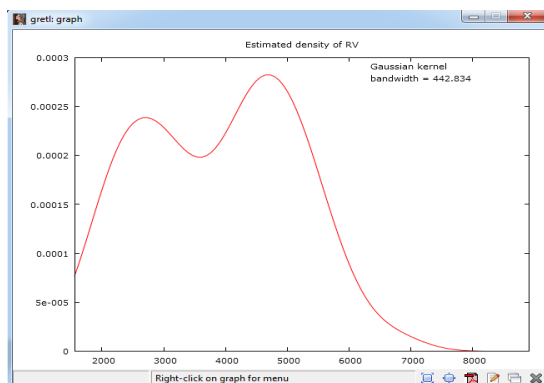
## Registered vehicles

a. Estimate the density functions.

To estimate the density function, highlight the variable of interest using the cursor and click

```
Variable --> Estimated density plot ...
```

and select the option *Gaussian kernel*.

RV                                                                RVSP

b. Frequency distribution and normality test.

RV



RVSP



c. Q-Q plots.

To obtain the Q-Q plot, highlight the variable of interest using the cursor, click

```
Variable --> Normal Q-Q plot ...
```

and mark use sample mean and variance for normal quantiles.



RV                                                          RVSP

d. Results:

- Both distributions have two modes.

- The Q-Q plots suggest that the distribution of the variable $RV$ is closer to the normal than the distribution of variable $RVSP$.
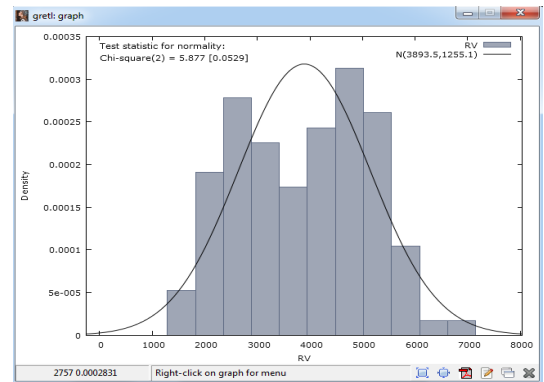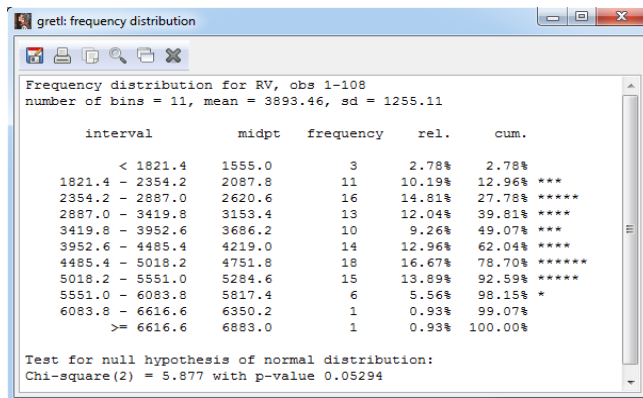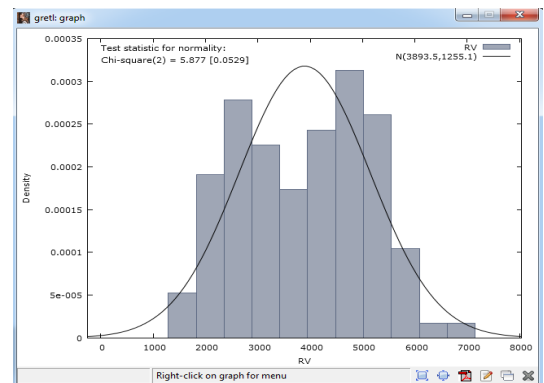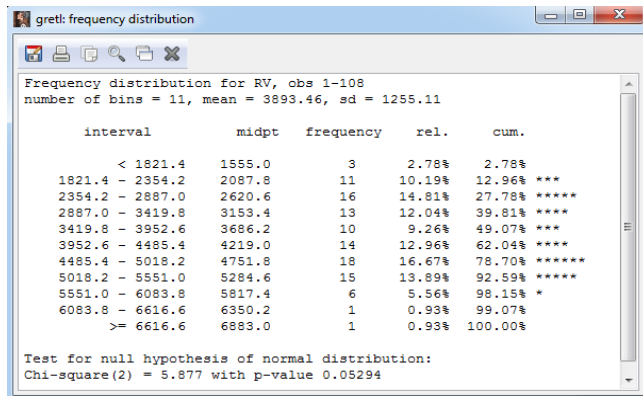
- The normality tests conclude that, at the $5\%$ significance level, the null hypothesis of normality is rejected for the variable $RVSP$ but it is not rejected for the variable $RV$.

Holiday cottages in Biscay

a. To analyse the distribution of the variable number of rooms, click

```
Variable --> Frequency distribution ...
```

and mark Show data only.



b. Frequency distribution of variable room price and normality tests.

<center>Normal distribution</center>



<center>Gamma distribution</center>

c. Frequency distribution of the qualitative variables.

LOC

```
gretl: frequency distribution

Frequency distribution for LOC, obs 1-75

        frequency    rel.     cum.

1          65       86.67%   86.67% ********************************
2          10       13.33%  100.00% ****
```

NP

```
gretl: frequency distribution

Frequency distribution for NP, obs 1-75

        frequency    rel.     cum.

1          34       45.33%   45.33% ****************
2          14       18.67%   64.00% ******
3          11       14.67%   78.67% *****
4          16       21.33%  100.00% *******
```
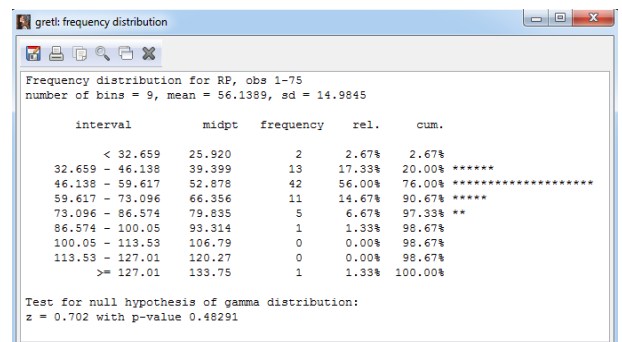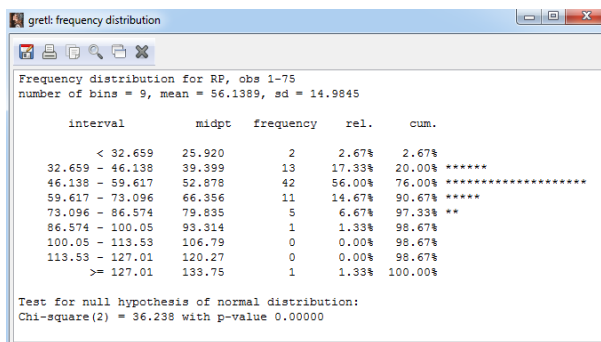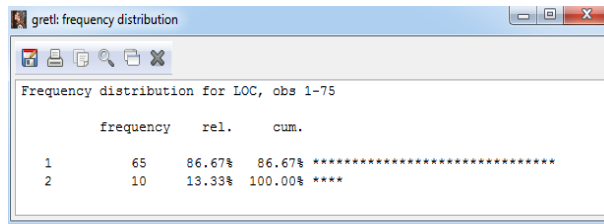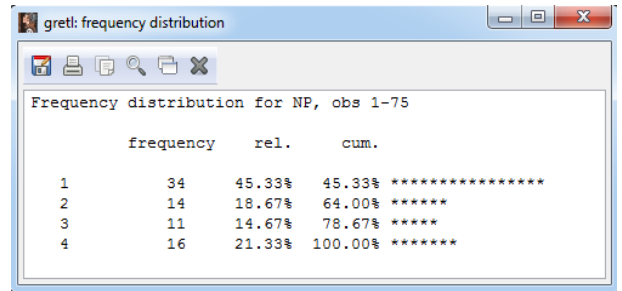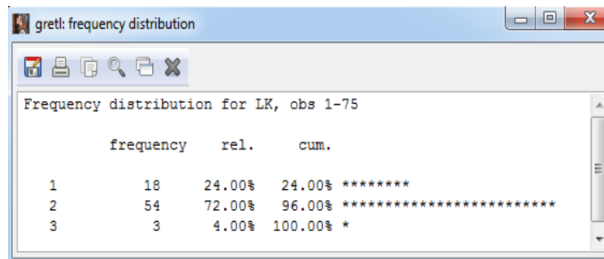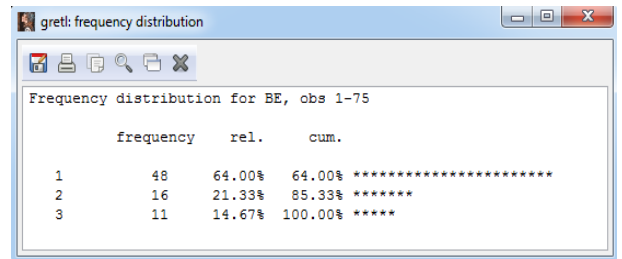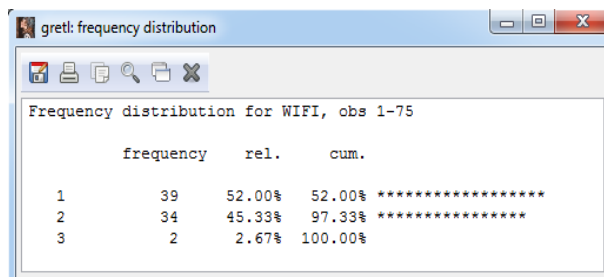
LK

```
gretl: frequency distribution

Frequency distribution for LK, obs 1-75

        frequency    rel.     cum.

1          18       24.00%   24.00% ********
2          54       72.00%   96.00% ************************
3           3        4.00%  100.00% *
```

BE

```
gretl: frequency distribution

Frequency distribution for BE, obs 1-75

        frequency    rel.     cum.

1          48       64.00%   64.00% ************************
2          16       21.33%   85.33% *******
3          11       14.67%  100.00% *****
```

WiFi

```
gretl: frequency distribution

Frequency distribution for WIFI, obs 1-75

        frequency    rel.     cum.

1          39       52.00%   52.00% ******************
2          34       45.33%   97.33% ****************
3           2        2.67%  100.00%
```

Q

```
gretl: frequency distribution

Frequency distribution for Q, obs 1-75

        frequency    rel.     cum.

1          73       97.33%   97.33% **************************************
2           2        2.67%  100.00%
```

C

```
gretl: frequency distribution

Frequency distribution for C, obs 1-75

        frequency    rel.     cum.

1          34       45.33%   45.33% ****************
2          41       54.67%  100.00% ********************
```

c.1. It refers to the category R = 3 of the variable BE. Therefore, 14.67 % of the holiday cottages are less than 1 km from a beach.

c.2. It refers to the category V = 1 of the variable LK. Therefore, 24 % of the holiday cottages are more than 20 km from a lake or reservoir.

c.3. It refers to the category P = 4 of the variable NP. Therefore, 21.33 % of the holiday cottages are inside a natural park.

c.4. It refers to the category R = 3 of the variable BE. Therefore, 14.67 % of the holiday cottages are less than 1 km from a beach.

It refers to the category R = 2 of the variable NP. Therefore, 18.67 % of the holiday cottages are less than 1 km from a natural park.

It refers to the category V = 1 of the variable BE. Therefore, 64 % of the holiday cottages are more than 20 km from a beach.

It refers to the category V = 1 of the variable NP. Therefore, 45.33 % of the holiday cottages are more than 20 km from a natural park.

c.5. It refers to the category A = 1 of the variable LOC. Therefore, 86.67 % of the holiday cottages are not in the the town center.

c.6. It refers to the categories G = 1 and S = 3 of the variable WiFi. Therefore, 54.67 % of the holiday cottages offer WiFi.

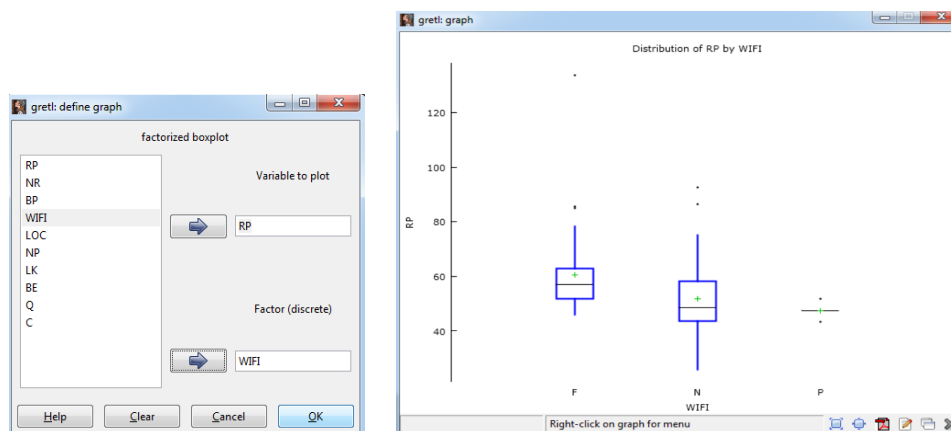It refers to the category G = 1 of the variable WiFi. Therefore, 52 % of the holiday cottages offer free WiFi.

c.7. It refers to the category N = 1 of the variable Q. Therefore, 97.33 % of the holiday cottages do not have the Tourism Quality Q certificate.

It refers to the category N = 2 of the variable C. Therefore, 54.67 % of the holiday cottages do not have the Commitment to Quality Tourism Award.

d. Boxplot of the variable average price of a room.

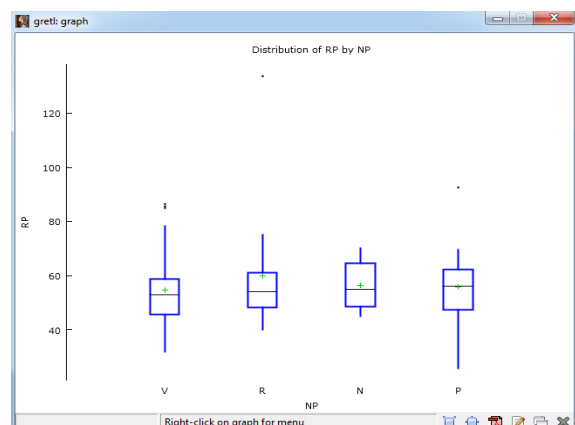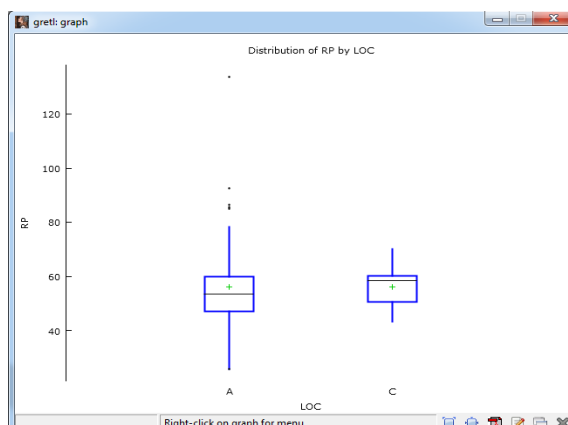$$\text{Variable} \; \texttt{-->} \; \text{Boxplot}$$

and mark <u>Factorized</u>. In the dialog box select the variable price of a room and the factor WiFi.
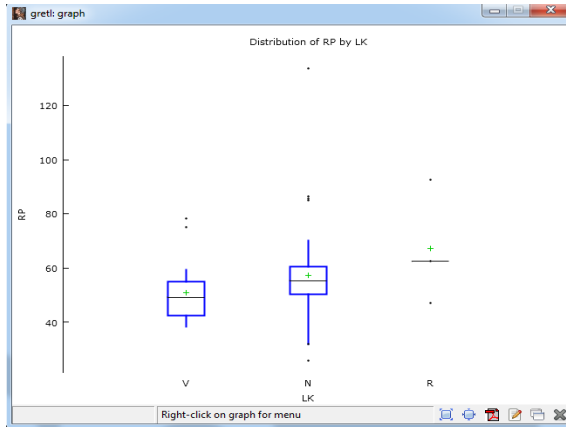


Using the rest of the qualitative variables as factors:

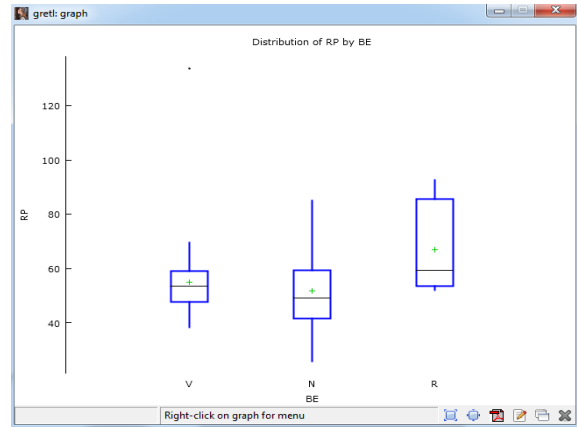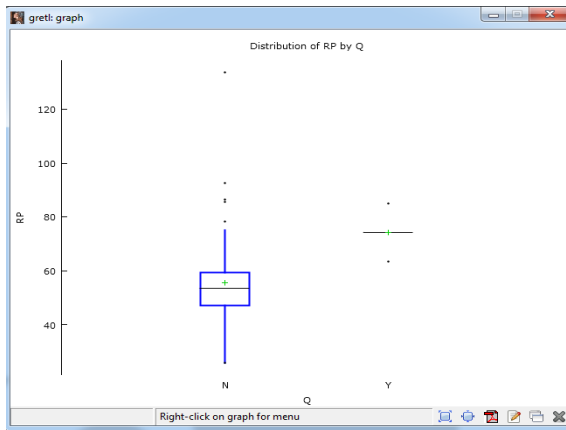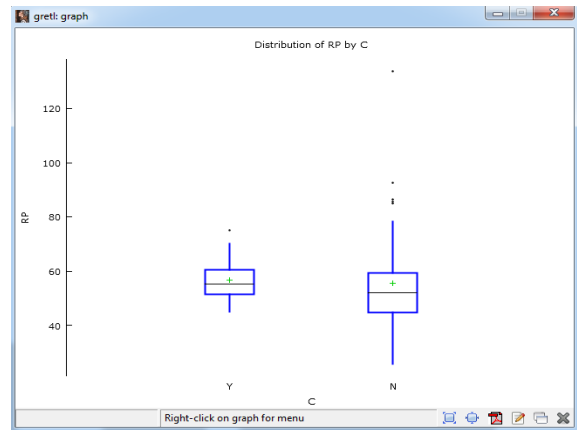LOC                                                                 NP
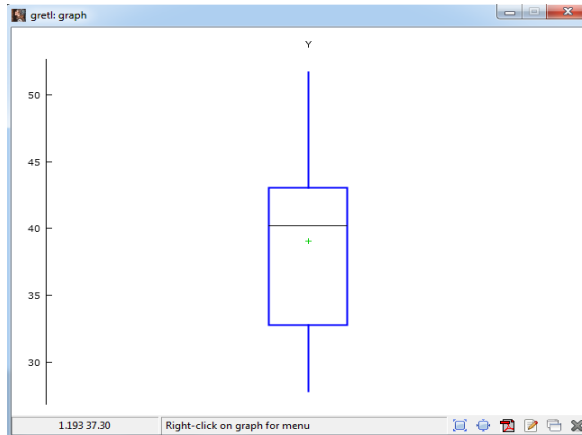
LK



BE



Q



C



e. Results.

- Given the results of the tests, the null hypothesis that the variable price follows a normal distribution is rejected at a 5 % significance level, but it is not rejected the null hypothesis that the variable price follows a gamma distribution.

- No.

- The frequency distribution of the variable LOC is concentrated in the category "far from the town center".

  The frequency distribution of the variable WiFi is concentrated in the categories "free WiFi access" and "No WiFi access".

- First, having the Tourism Quality Q certificate, followed by WiFi access and proximity to a lake or a beach.

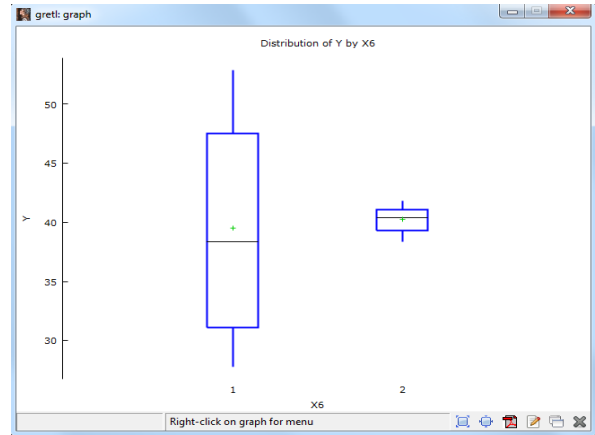## Chicken consumption

a. and b. Boxplots.

<table>
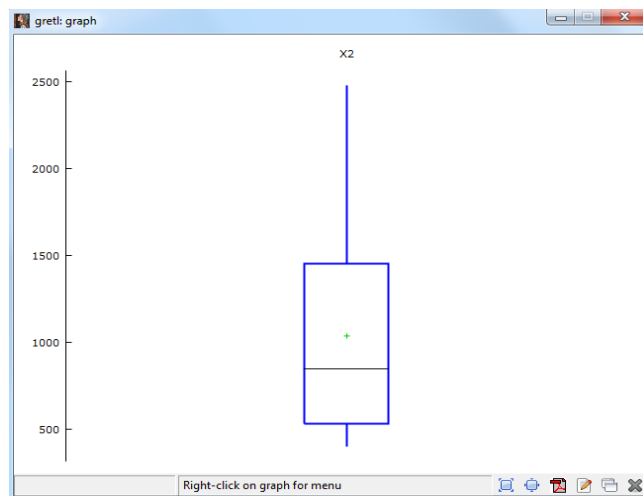<tr><td align="center">Simple boxplot</td><td align="center">Factorized boxplot (X6)</td></tr>
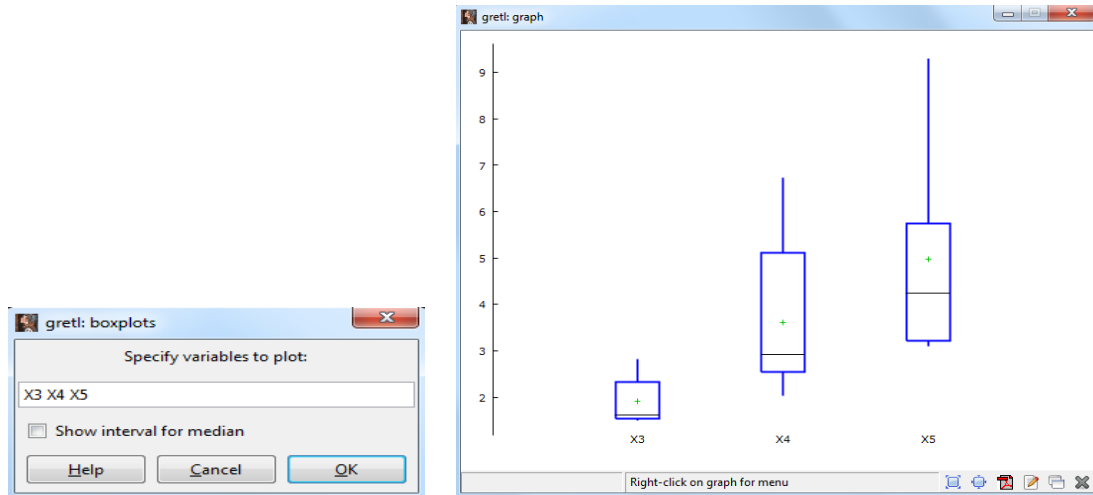</table>


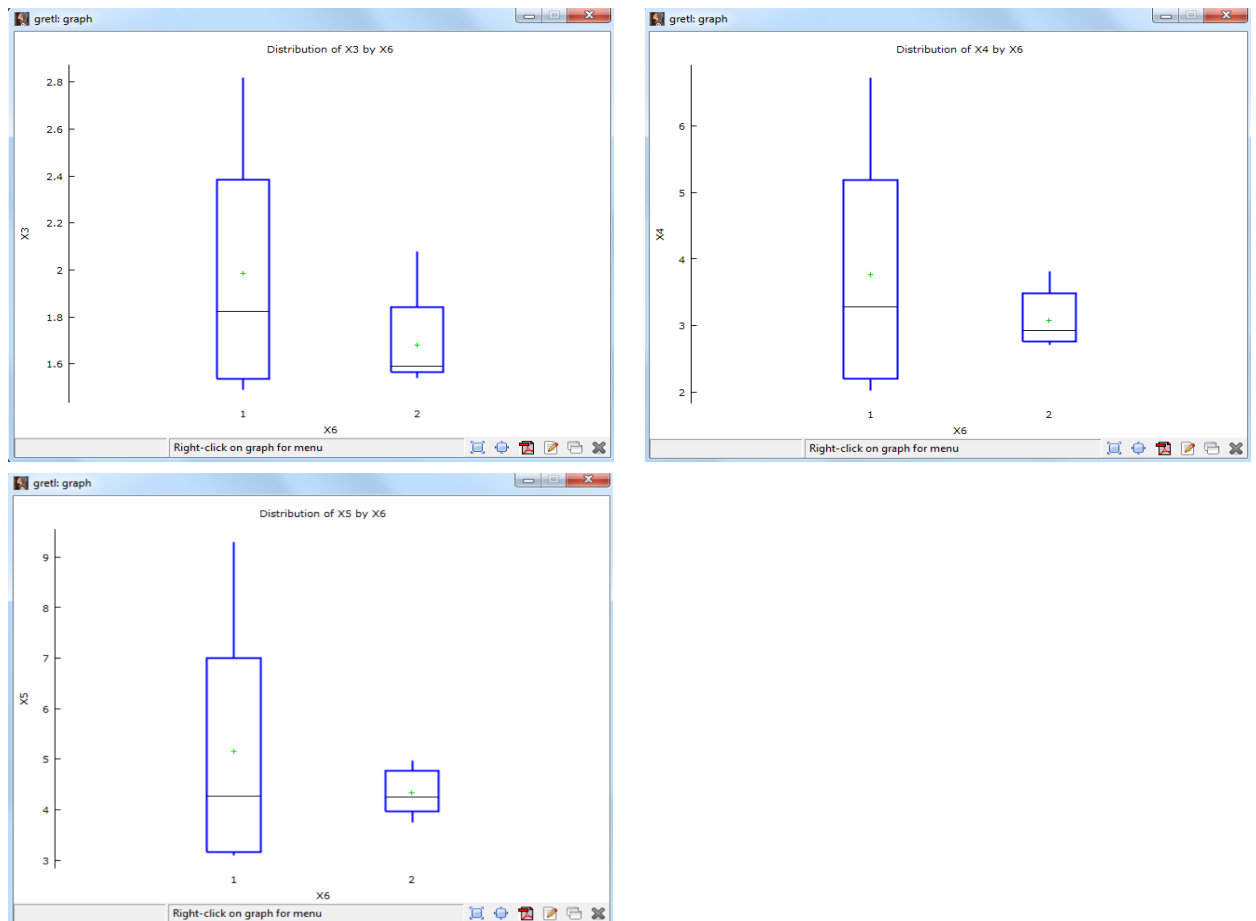
c. Boxplot of variable income.

```
Variable --> Boxplot
```

and mark Simple boxplot.



d. Joint boxplot of all the price variables.

```
View --> Graphs specified vars    --> Boxplots ...
```

Then, select the price variables.

e. Factorized boxplot for each price variable.



f. Results.

- The variability observed in the variable consumption of chicken comes mostly from the period of avian flu epidemic. Note that the variability observed in the years without epidemic is very small.

- The dispersion in the distribution of the variable chicken price is much smaller than the dispersion in the distributions of the beef price and of the pork price.

- The biggest difference in the price of a product due to the avian flue epidemic is observed in the price of chicken.