

MODS202: Project's Report - Part 1

Due on Sunday, November 22, 2020

Professor Patrick Waelbroeck

Daniel Deutsch and Kevin Kühl

Answers of Part 1

Question 1

The data was loaded and cleaned. In that process, only entries with a positive value for wage were kept. A brief visualization of the final loaded data can be seen below, where some of the columns do not appear for size constraints.

	inlf	hours	kidslt6	kidsgt6	age	educ	wage	repwage	hushrs	husage	...	faminc	mtr	motheduc	fatheaduc	unem	city	exper	nwifeinc
0	1.0	1610.0	1.0	0.0	32.0	12.0	3.3540	2.65	2708.0	34.0	...	16310.0	0.7215	12.0	7.0	5.0	0.0	14.0	10.910060
1	1.0	1656.0	0.0	2.0	30.0	12.0	1.3889	2.65	2310.0	30.0	...	21800.0	0.6615	7.0	7.0	11.0	1.0	5.0	19.499980
2	1.0	1980.0	1.0	3.0	35.0	12.0	4.5455	4.04	3072.0	40.0	...	21040.0	0.6915	12.0	7.0	5.0	0.0	15.0	12.039910
3	1.0	456.0	0.0	3.0	34.0	12.0	1.0965	3.25	1920.0	53.0	...	7300.0	0.7815	7.0	7.0	5.0	0.0	6.0	6.799996
4	1.0	1568.0	1.0	2.0	31.0	14.0	4.5918	3.60	2000.0	32.0	...	27300.0	0.6215	12.0	14.0	9.5	1.0	7.0	20.100060
...
423	1.0	680.0	0.0	5.0	36.0	10.0	2.3118	0.00	3430.0	43.0	...	19772.0	0.7215	7.0	7.0	7.5	0.0	2.0	18.199980
424	1.0	2450.0	0.0	1.0	40.0	12.0	5.3061	6.50	2008.0	40.0	...	35641.0	0.6215	7.0	7.0	5.0	1.0	21.0	22.641060
425	1.0	2144.0	0.0	2.0	43.0	13.0	5.8675	0.00	2140.0	43.0	...	34220.0	0.5815	7.0	7.0	7.5	1.0	22.0	21.640080
426	1.0	1760.0	0.0	1.0	33.0	12.0	3.4091	3.21	3380.0	34.0	...	30000.0	0.5815	12.0	16.0	11.0	1.0	14.0	23.999980
427	1.0	490.0	0.0	1.0	30.0	12.0	4.0816	2.46	2430.0	33.0	...	18000.0	0.6915	12.0	12.0	7.5	1.0	7.0	16.000020

428 rows × 22 columns

After the cleaning and formatting process, the data was composed by 428 entries, each containing 22 variables (represented by the columns).

Question 2

First, the descriptive statistics of wage, age and education for the whole set of females.

	wage	age	educ
count	428.000000	428.000000	428.000000
mean	4.177682	41.971963	12.658879
std	3.310282	7.721084	2.285376
min	0.128200	30.000000	5.000000
25%	2.262600	35.000000	12.000000
50%	3.481900	42.000000	12.000000
75%	4.970750	47.250000	14.000000
max	25.000000	60.000000	17.000000

Note worthy points

- The mean wage is 4.18 (hourly)
- The mean age is 41.98 years
- The mean number of years of education is 12.66

Selecting now only the entries in which the wage of the husband is greater than the median of wage of the whole set of husbands.

	wage	age	educ
count	214.000000	214.000000	214.000000
mean	4.896822	42.275701	13.242991
std	4.041606	7.388843	2.359045
min	0.161600	30.000000	5.000000
25%	2.513850	36.000000	12.000000
50%	3.846400	43.000000	12.000000
75%	5.854125	48.000000	16.000000
max	25.000000	59.000000	17.000000

Note worthy points

- The mean wage is 4.89 (hourly)
- The mean age is 42.27 years
- The mean number of years of education is 13.24

Finally, we select the entries in which the wage of the husband is less than the median of wage of the whole set of husbands.

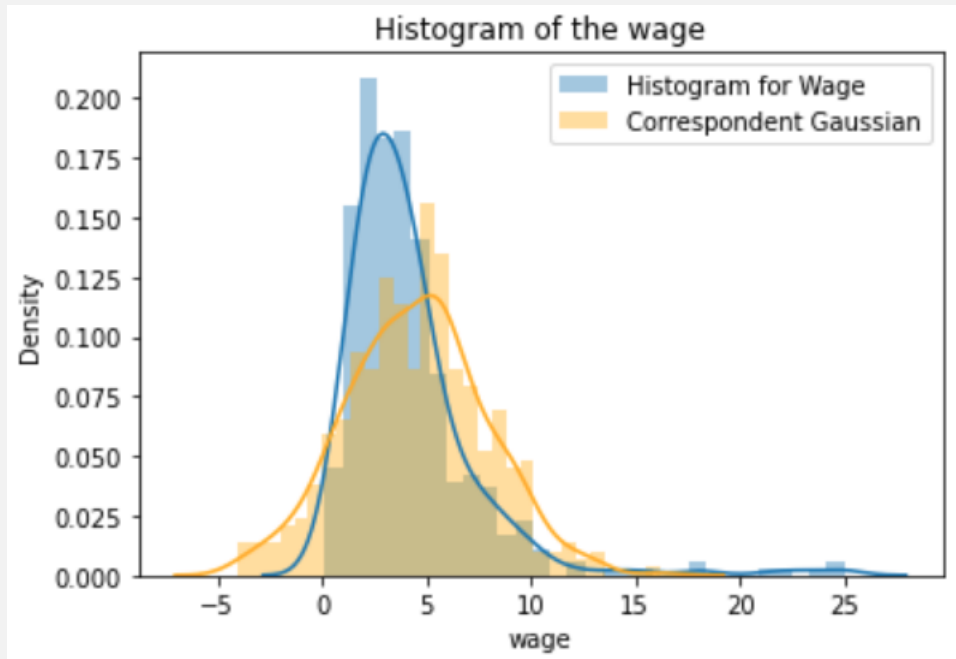
	wage	age	educ
count	214.000000	214.000000	214.000000
mean	3.458541	41.668224	12.074766
std	2.143274	8.045482	2.054200
min	0.128200	30.000000	6.000000
25%	2.117275	35.000000	12.000000
50%	2.971800	41.000000	12.000000
75%	4.393800	47.000000	12.000000
max	18.267000	60.000000	17.000000

Note worthy points

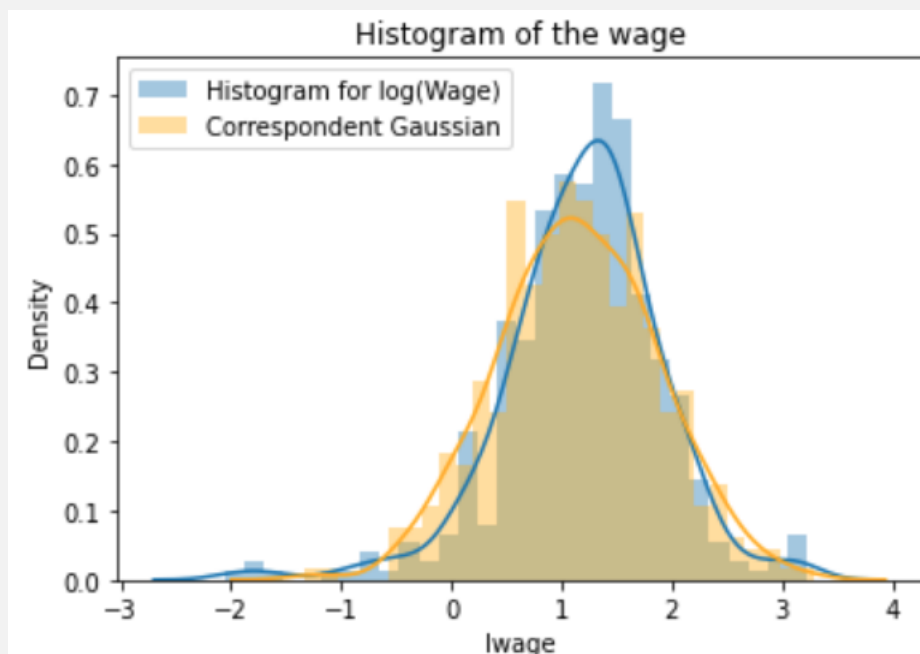
- The mean wage is 3.46 (hourly)
- The mean age is 41.67 years
- The mean number of years of education is 12.07

Question 3

We plot the histogram of the variable wage together with the ideal gaussian distribution which would represent it based on the mean and standard deviation of the data



After calculating the logarithm of the variable wage, we get the following histogram



It can be noted that the histogram of the $\log(\text{wage})$ fits the gaussian distribution in a better way than the histogram of the wage.

Question 4

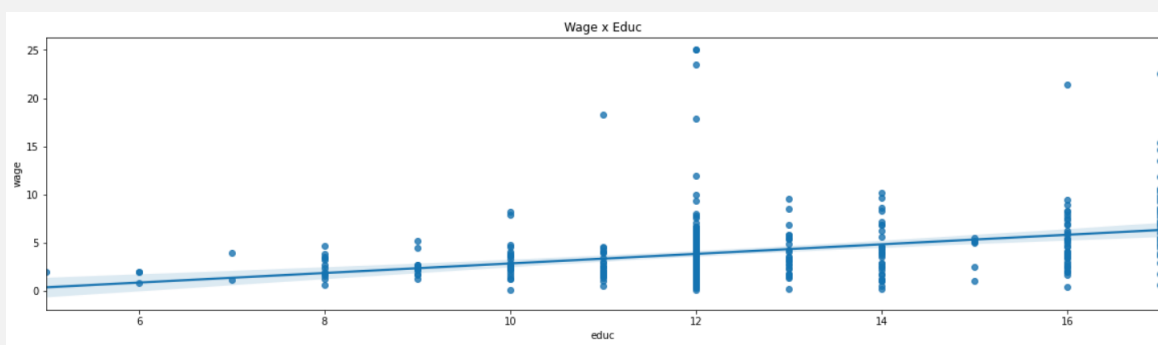
If we calculate the correlation between the variables motheduc and fatheduc, we obtain the following result

	motheduc	fatheduc
motheduc	1.000000	0.554063
fatheduc	0.554063	1.000000

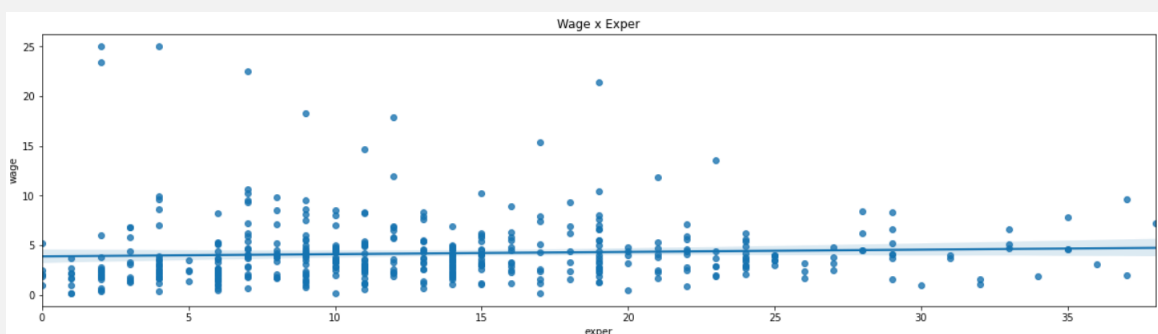
Therefore, the correlation between the two variables is 0.554063, which is high. This can be understood as usually people marry in within their socioeconomic class, which, in general, achieves similar degree of education. A multicollinearity problem will rise if one uses both variables as explanatory variables for a predictive model.

Question 5

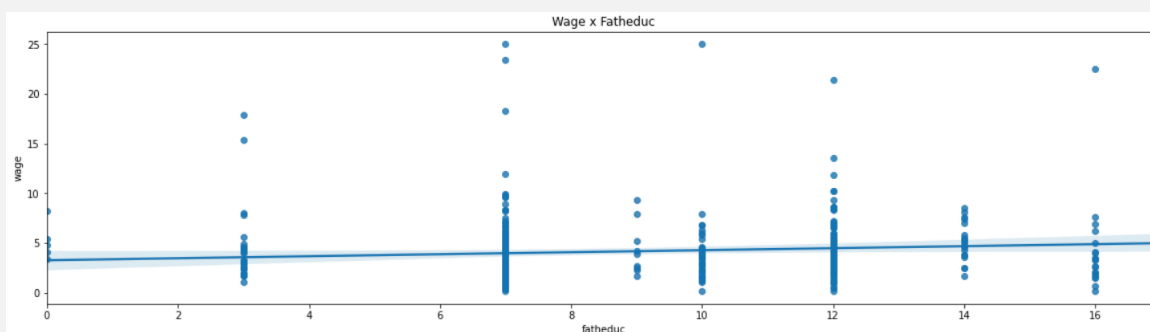
The plot for the variables Wage and Educ is given below



Now, for the variables Wage and Exper



Finally, for the variables Wage and Fatheduc



First plot indicates a clear relation between Wage and Education. In this case, the higher the number the years of education, the higher the wage tends to be.

In the second plot, the regression between the two variables show a very lightly inclined line, which would suggest a soft positive relation between wage and the number of years of experience.

For the last plot, there is also a small inclination, which would suggest that the higher the number of years of education a father has, the higher the wage of his daughter tends to be. The slope, however, like last plot, is very small.

This is not a case of the effect "toute chose étant égale par ailleurs" because we are not keeping the other variables constant when analysing the two plotted variables. For example, in the first plot, we are not considering the number of kids a certain individual, which can vary for each point.

Question 6

The fundamental hypothesis for an unbiased estimator is that the non observed variable have a null mean and that the conditional mean of the non observed variable given the data is equal to the unconditional mean (which is null). In summary, we can write

$$E(u|x) = E(u) = 0$$

The "biais de variable omise", or omitted variable bias is the circumstance in which relevant variables are left out of the regression model. For instance, these variables are independent from the others used and can carry great information about the overall data.

Question 7

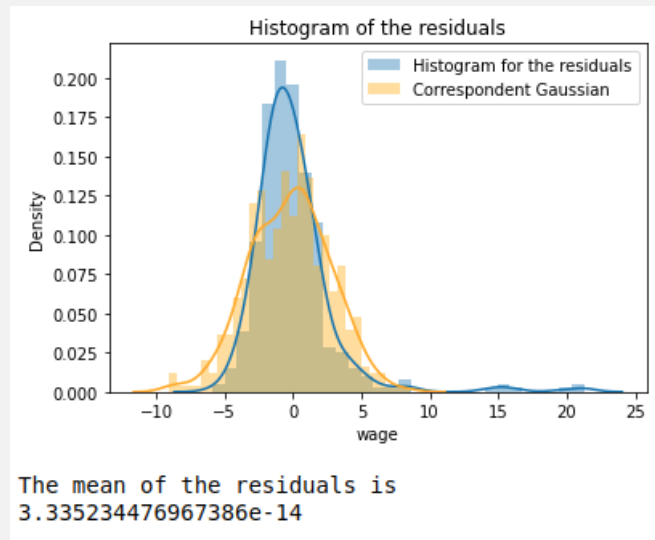
After constructing the regression, using the dependent variable wage and explanatory variables city, educ, exper, nwifeinc, kidslt6, kidsgt6, we get the following result

$$wage = c + \beta_1 city + \beta_2 educ + \beta_3 exper + \beta_4 nwifeinc + \beta_5 kidslt6 + \beta_6 kidsgt6$$

With

$$\begin{pmatrix} c \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \beta_6 \end{pmatrix} = \begin{pmatrix} -2.403454 \\ 0.369752 \\ 0.460048 \\ 0.023820 \\ 0.015245 \\ 0.036173 \\ -0.061891 \end{pmatrix}$$

The histogram of residuals for the above regression is given below.



The residuals have some components that make them not very symmetric around 0. Also, the variance of these residuals is quite high (9.700).

Question 8

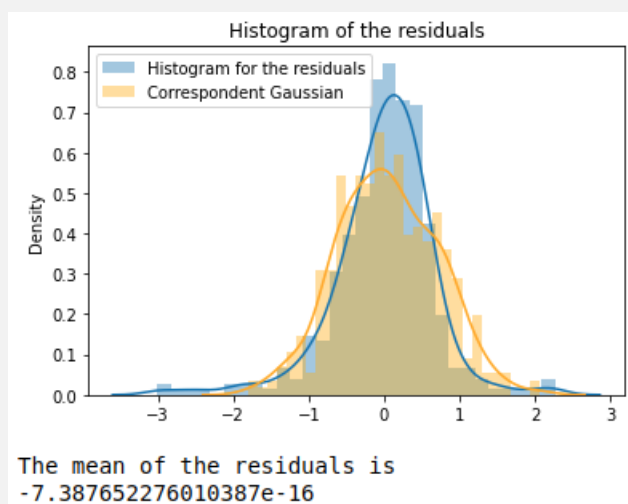
After constructing the regression, using the dependent variable $\ln wage$ and explanatory variables $city$, $educ$, $exper$, $nwifeinc$, $kidslt6$, $kidsgt6$, we get the following result

$$\ln(wage) = c + \beta_1 city + \beta_2 educ + \beta_3 exper + \beta_4 nwifeinc + \beta_5 kidslt6 + \beta_6 kidsgt6$$

With

$$\begin{pmatrix} c \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \beta_6 \end{pmatrix} = \begin{pmatrix} -0.398975 \\ 0.035268 \\ 0.102248 \\ 0.015488 \\ 0.004883 \\ -0.045303 \\ -0.011704 \end{pmatrix}$$

The histogram of residuals for the above regression is given below.



We can observe that the residuals are better distributed around zero. They better fit a normal distribution associated to the given residuals mean and standard deviation. The variance of the residuals have significantly decreased (now it is 0.4479). Also the mean of the residuals, compared to last case, is more close to zero.

Question 9

In this hypothesis test, we make

$$H_0 : \beta_4 = \text{zero}$$

$$H_1 : \beta_4 \neq 0$$

Considering the regression made in question 8, we get the following results for the t-student test for the variables, in specific for nwifeinc.

OLS Regression Results						
Dep. Variable:	lwage	R-squared:	0.156			
Model:	OLS	Adj. R-squared:	0.144			
Method:	Least Squares	F-statistic:	12.92			
Date:	Sat, 21 Nov 2020	Prob (F-statistic):	2.00e-13			
Time:	20:14:13	Log-Likelihood:	-431.92			
No. Observations:	428	AIC:	877.8			
Df Residuals:	421	BIC:	906.3			
Df Model:	6					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	-0.3990	0.207	-1.927	0.055	-0.806	0.008
x1	0.0353	0.070	0.503	0.616	-0.103	0.173
x2	0.1022	0.015	6.771	0.000	0.073	0.132
x3	0.0155	0.004	3.452	0.001	0.007	0.024
x4	0.0049	0.003	1.466	0.143	-0.002	0.011
x5	-0.0453	0.085	-0.531	0.596	-0.213	0.122
x6	-0.0117	0.027	-0.434	0.664	-0.065	0.041
Omnibus:	79.542	Durbin-Watson:	1.979			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	287.193			
Skew:	-0.795	Prob(JB):	4.33e-63			
Kurtosis:	6.685	Cond. No.	178.			

Observing the above results, we can conclude the results of the hypothesis testing looking at the p-value ($P > |t|$). This is already the two-tailed test p-value. **For the t-student test for variable nwifeinc, two-tailed test p-value is 0.143.**

- For alpha = 1%, we consider that we must have a p-value < 0.01 to reject the null hypothesis (coefficient of nwifeinc is zero). **We do not reject it.**

- For alpha = 5%, we consider that we must have a p-value < 0.05 to reject the null hypothesis (coefficient of nwifeinc is zero). **We do not reject it.**

- For alpha = 10%, we consider that we must have a p-value < 0.10 to reject the null hypothesis (coefficient of nwifeinc is zero). **We do not reject it.**

Therefore, we can conclude that the variable nwifeinc is not significant with test significance level of 1%. This is the same as saying that it is not relevant for explaining the variable log(wage).

Question 10

In this hypothesis test, we make

$$H_0 : \beta_4 = 0.01$$

$$H_1 : \beta_4 \neq 0.01$$

We execute the test by translating the beta values before dividing it by the standard deviation. We obtain.

```
T-test results
[-1.97524409  0.36005271  6.10838827  1.22305946 -1.53638899 -0.64827498
 -0.80549245]

In specific for the nwifeinc variable's coefficient
-1.536388955562578

T-test p-value for the nwifeinc variable's coefficient
0.12519597367688684
```

Observing the above results, we can conclude the results of the hypothesis testing looking at the p-value (already considering a two-tailed test). **For the test of nwifeinc == 0.01, two-tailed test p-value is 0.1252.**

- For alpha = 5%, we consider that we must have a p-value < 0.05 to reject the null hypothesis (that the coefficient of nwifeinc is 0.01). **We do not reject it.**

Therefore, we can conclude that the coefficient of the variable nwifeinc is equal to 0.01 with test significance level of 5%.

Question 11

In this hypothesis test, we make

$$H_0 : \beta_4 = 0.01 \quad \text{and} \quad \beta_1 = 0.05$$

$$H_1 : \beta_4 \neq 0.01 \quad \text{or} \quad \beta_1 \neq 0.05$$

The Fisher statistic is calculated by

$$F = \frac{\frac{SSR_{constrained} - SSR_{unconstrained}}{q}}{\frac{SSR_{unconstrained}}{(n-k-1)}}$$

In our particular case, we have $q=2$.

Applying it, we get as result

```
SSR0 is equal to: 188.58998019263944   Unconstrained model
SSR1 is equal to: 189.78788085217226   Constrained model Beta_4 = 0.01 and Beta_1 = 0.05
The Fisher test result: 1.3370704454928417
The p-value: 0.2637267136252716
```

For the test, we get as result a p-value of 0.2637

- For $\alpha = 5\%$, we consider that we must have a p-value < 0.05 to reject the null hypothesis (that the coefficient of `nwifeinc` is 0.01 and coefficient of `city` is 0.05). **We do not reject it.**

Therefore, we can conclude that the coefficient of the variable `nwifeinc` is equal to 0.01 and the coefficient of the variable `city` is equal to 0.05 with test significance level of 5%.

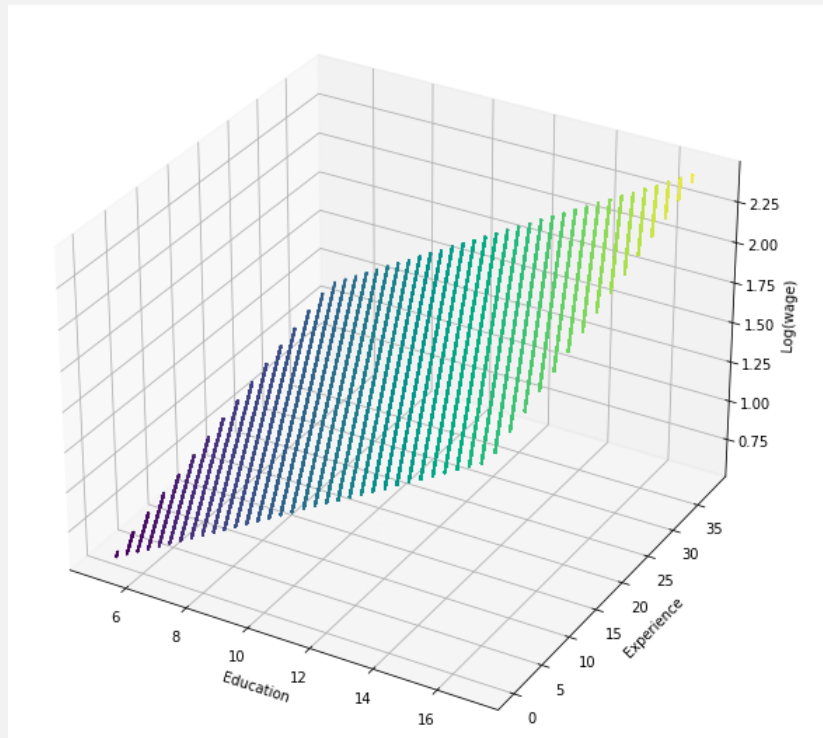
Question 12

Let's consider the regression of $\log(wage)$ considering as explanatory variables the number of years of education and the number of years of experience. Where, in the following, x_1 represents the education variable and x_2 represents the experience variable.

OLS Regression Results						
Dep. Variable:	lwage	R-squared:	0.148			
Model:	OLS	Adj. R-squared:	0.144			
Method:	Least Squares	F-statistic:	37.02			
Date:	Sun, 22 Nov 2020	Prob (F-statistic):	1.51e-15			
Time:	12:14:05	Log-Likelihood:	-433.74			
No. Observations:	428	AIC:	873.5			
Df Residuals:	425	BIC:	885.6			
Df Model:	2					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	-0.4002	0.190	-2.102	0.036	-0.774	-0.026
x1	0.1095	0.014	7.728	0.000	0.082	0.137
x2	0.0157	0.004	3.900	0.000	0.008	0.024
Omnibus:	81.122	Durbin-Watson:	1.981			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	296.773			
Skew:	-0.807	Prob(JB):	3.60e-65			
Kurtosis:	6.746	Cond. No.	113.			

We can plot the regression (1 dependent variable and 2 explanatory variables) in a 3 dimen-

sional plot, where the z-axis is represented by the dependent variable.



It is clear that the higher the years of education and the years of experience, the higher the wage.

Also, if we fix one value for one of the explanatory variables, we can see that the regression is positive related with the other. For example, fixing a value for the number of years of education, we observe that the higher the number of years of experience the higher the wage.

Question 13

Initial model

$$y = \log(wage) = c + \beta_1 city + \beta_2 educ + \beta_3 exper + \beta_4 nwifeinc + \beta_5 kidslt6 + \beta_6 kidsgt6$$

Defining

$$\theta = \beta_6 - \beta_5$$

We can write

$$\theta + \beta_5 = \beta_6$$

We get the model

$$\log(wage) = c + \beta_1 city + \beta_2 educ + \beta_3 exper + \beta_4 nwifeinc + \beta_5 kidslt6 + (\theta + \beta_5) kidsgt6$$

Which is, finally written as

$$\log(wage) = c + \beta_1 city + \beta_2 educ + \beta_3 exper + \beta_4 nwifeinc + \beta_5 (kidslt6 + kidsgt6) + \theta kidsgt6$$

Now, we test it for the significance of θ . Therefore we write

$$H_0 : \theta = 0$$

$$H_1 : \theta \neq 0$$

OLS Regression Results						
Dep. Variable:	lwage	R-squared:	0.156			
Model:	OLS	Adj. R-squared:	0.144			
Method:	Least Squares	F-statistic:	12.92			
Date:	Sat, 21 Nov 2020	Prob (F-statistic):	2.00e-13			
Time:	20:14:13	Log-Likelihood:	-431.92			
No. Observations:	428	AIC:	877.8			
Df Residuals:	421	BIC:	906.3			
Df Model:	6					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	-0.3990	0.207	-1.927	0.055	-0.806	0.008
x1	0.0353	0.070	0.503	0.616	-0.103	0.173
x2	0.1022	0.015	6.771	0.000	0.073	0.132
x3	0.0155	0.004	3.452	0.001	0.007	0.024
x4	0.0049	0.003	1.466	0.143	-0.002	0.011
x5	-0.0453	0.085	-0.531	0.596	-0.213	0.122
x6	0.0336	0.090	0.372	0.710	-0.144	0.211
Omnibus:	79.542	Durbin-Watson:	1.979			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	287.193			
Skew:	-0.795	Prob(JB):	4.33e-63			
Kurtosis:	6.685	Cond. No.	178.			

Therefore, **the p-value for the two-tailed significance test of θ is 0.710.**

- For alpha = 5%, we consider that we must have a p-value < 0.05 to reject the null hypothesis (coefficient of kidsgt6 and coefficient of kidslt6 are equal). **We do not reject it.**

Therefore, we can conclude that $\theta = 0$ with significance level of 5%. This means that the coefficients associated to kidsgt6 and kidslt6 are equal, with significance level of 5%.

Intepreting this, we get that the effect of children in the wage of a female is the same, disregarding the fact the children are young (less than 6 years) or old ($6 < \text{age} < 18$).

Question 14

In this case we make

H_0 : The data has homoscedasticity.

H_1 : The data has heteroscedasticity.

We propose the test that suppose a linear relation between the squared error term (residuals) and the used variables.

So, we suppose

$$u^2 = \delta_0 + \delta_1 x_1 + \dots \delta_k x_k + v$$

We test, therefore

$$H_0 : \delta_1 = \delta_2 = \dots = \delta_k = 0$$

As a result of the regression in u^2 we get

OLS Regression Results						
Dep. Variable:	y	R-squared:	0.022			
Model:	OLS	Adj. R-squared:	0.008			
Method:	Least Squares	F-statistic:	1.593			
Date:	Sat, 21 Nov 2020	Prob (F-statistic):	0.148			
Time:	21:05:11	Log-Likelihood:	-2207.4			
No. Observations:	428	AIC:	4429.			
Df Residuals:	421	BIC:	4457.			
Df Model:	6					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	1.4856	13.111	0.113	0.910	-24.285	27.256
x1	5.9644	4.444	1.342	0.180	-2.770	14.699
x2	0.8077	0.956	0.845	0.399	-1.072	2.687
x3	-0.5341	0.284	-1.880	0.061	-1.093	0.024
x4	0.0435	0.211	0.206	0.837	-0.371	0.458
x5	4.9573	5.402	0.918	0.359	-5.661	15.575
x6	-0.4018	1.706	-0.236	0.814	-3.756	2.952
Omnibus:	638.793	Durbin-Watson:	2.029			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	96122.227			
Skew:	8.127	Prob(JB):	0.00			
Kurtosis:	74.595	Cond. No.	178.			

As explained in the project's instructions, all tests must be done considering a 5% significance level. We obtained a p-value of 0.148 for the Linearity test for heterocedascity.

Therefore, within the required significance, we cannot reject H_0 , meaning that we can not conclude that there is heterocedascity in the data (considering the model of question 7).

Question 15

To test the change of structure on the regression made on question 8, we propose the use of the Chow test. Under this test, we make

H_0 : There is a change in structure

H_1 : There is not a change in structure

The test statistic is calculated as follow (for two groups)

$$F_{Chow} = \frac{SSR_0 - (SSR_1 + SSR_2)}{SSR_1 + SSR_2} \times \frac{n_1 + n_2 - 2k}{k}$$

With n being the number of observations in a given group, and k is the number of parameters we are estimating.

- First, with the following groups
 - Entire set
 - Women with age greater or equal than 43 years old
 - Women with less than 43 years old

For the proposed test, we get the following results

```
SSR0: 188.58998019263944
SSR1: 80.40365115321053
SSR2: 104.48165074506036
n: 428
n1: 211
n2: 217
k: 7

Chow Test p-value: 0.3099734135726031
Fisher: 1.1850874941083283
```

Therefore, with a p-value of 0.3099, **we accept the null hypothesis (with a significance level of 5%) that there exist a change in structure for between the group of women with more than 43 years and the group of women with less than 43 years.**

- for the second part, we propose the following groups
 - Entire set
 - Women with less than 30 years old (\leq)
 - Women with more than 30 years old and less than 43 years old ($>$ and $<$)
 - Women with more than 43 years old (\geq)

For this test, we get the following results

```
SSR0: 188.58998019263944
SSR1: 3.387215073161979
SSR2: 100.03565630367271
SSR3: 80.40365115321053
n: 428
n1: 19
n2: 198
n3: 211
k: 7

Chow Test p-value: 0.0961451078434979
Fisher: 1.5325563045387653
```

Therefore, with a p-value of 0.0961, **we accept the null hypothesis (with a significance level of 5%) that there exist a change in structure for between the proposed groups.**

Question 16

In order to do the regression we write the variables

- more43: 1 if woman's age is ≥ 43 ; 0 otherwise
- between3043: 1 if woman's age is >30 but <43 ; 0 otherwise
- less30: 1 if woman's age is ≤ 30 ; 0 otherwise

We cannot add those three variables to the model as

$$more43 + between3043 + less30 = 1$$

We would have a multicollinearity problem.

Therefore, using as the base case the women that have less than 30 years old (\leq) we build the regression adding more43 and between3043.

The regression yields (coefficients for the binary variables in the last two lines).

OLS Regression Results						
Dep. Variable:	lwage	R-squared:	0.161			
Model:	OLS	Adj. R-squared:	0.145			
Method:	Least Squares	F-statistic:	10.07			
Date:	Sun, 22 Nov 2020	Prob (F-statistic):	7.08e-13			
Time:	14:22:57	Log-Likelihood:	-430.46			
No. Observations:	428	AIC:	878.9			
Df Residuals:	419	BIC:	915.5			
Df Model:	8					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	-0.2210	0.248	-0.893	0.373	-0.708	0.266
x1	0.0475	0.071	0.672	0.502	-0.092	0.187
x2	0.1008	0.015	6.653	0.000	0.071	0.131
x3	0.0179	0.005	3.785	0.000	0.009	0.027
x4	0.0058	0.003	1.712	0.088	-0.001	0.012
x5	-0.0809	0.088	-0.920	0.358	-0.254	0.092
x6	-0.0183	0.029	-0.642	0.521	-0.074	0.038
x7	-0.2558	0.169	-1.513	0.131	-0.588	0.077
x8	-0.1618	0.164	-0.986	0.325	-0.484	0.161
Omnibus:	77.107	Durbin-Watson:	1.996			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	282.164			
Skew:	-0.764	Prob(JB):	5.36e-62			
Kurtosis:	6.673	Cond. No.	254.			

This means

$$lwage = c + \beta_1 city + \beta_2 educ + \beta_3 exper + \beta_4 nwifeinc + \beta_5 kidslt6 + \beta_6 kidsgt6 + \beta_7 more43 + \beta_8 between3043$$

We perform the Fisher test with hypothesis

$$H_0 : \beta_7 = 0 \quad \text{and} \quad \beta_8 = 0$$

$$H_1 : \beta_7 \neq 0 \quad \text{or} \quad \beta_8 \neq 0$$

As the results for the test we obtain

```
Value of q: 2
n: 428 / k: 7
SSR0 is equal to: 187.30681260001555
SSR1 is equal to: 188.58998019263944
The Fisher test result: 1.442055280840906
The p-value: 0.23760676040907033
```

Finally, with a p-value of 0.2376, we do not reject the null hypothesis with test significance level of 5%. Therefore, both variables are not significant.