

MODS PROJECT

Imports

```
In [127]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import statsmodels.api as sm
from scipy.stats import t
from scipy.stats import spearmanr
from scipy.stats import norm

In [128]: df = pd.read_csv("mod2.txt", delimiter=True, header=None)
```

Naming columns

```
In [129]: df = df.rename(columns={0: 'laborforce'
, 1: 'hours worked'
, 2: 'kids below 6'
, 3: 'kids between 6 and 18'
, 4: 'age women'
, 5: 'educ years of schooling'
, 6: 'wage estimated'
, 7: 'husbandage'
, 8: 'hours worked by husband'
, 9: 'husband age'
, 10: 'husband hour/w wage'
, 11: 'husband tax rate facing women'
, 12: 'family income'
, 13: 'marginal tax rate facing women'
, 14: 'mothers year of schooling'
, 15: 'fathers year of schooling'
, 16: 'rate in country of resit'
, 17: 'live in SMSA'
, 18: 'marginal tax rate facing woman'
, 19: 'mfefinc'
, 20: 'log(wage)'
, 21: 'experq'})

Out[129]:
```

```

752      0      0      0      3      39      9      0.00      3120      48      28363      0.6915      7      7      11.0      1      12      28.363000
753 rows x 22 columns

Question 1
Let's transform " " into real values

In [180]: s.df = df[s['wage estimated']!=0].copy()
empty_rows = s[s['wage' == ''].index
df = df.drop(empty_rows)
df=df.copy()

In [181]: print("We removed "+len(df0)-len(df)," where the wage data was null or negative")

We removed 0 where the wage data was null or negative

Question 2
For all women:

In [182]: df['wage estimated'] = pd.to_numeric(df['wage estimated'])
print("Statistics of women salaries whose husband earn more the median are:")
df[['wage estimated', 'age women', 'educ years of schooling']].describe()

The statistics of women salaries whose husband earn more the median are:

```

Question 1

```
In [128]: s = df[['wage estimated'], df0.copy()]
empty_rows = s[s=="."]
df = df[empty_rows.index]
df = df[empty_rows.index]
```

```
In [130]: print("We removed ", len(df0)-len(df), " where the wage data was null or negative")

We removed 8 where the wage data was null or negative
```

Question 2

For all women:

```
In [132]: df["wage estimated"] = pd.to_numeric(df["wage estimated"])
df["wage estimated"] = df["wage estimated"].replace(0, np.nan)
df["wage estimated"] = df["wage estimated"].replace(0, np.nan)
```

The statistics of women salaries whose husband earn more the median are:

| | wage estimated | age women | educ years of schooling |
|-------|----------------|------------|-------------------------|
| count | 428.000000 | 428.000000 | 428.000000 |
| mean | 4.177682 | 41.971963 | 12.658979 |
| std | 3.310282 | 7.722084 | 2.283376 |
| min | 0.000000 | 30.000000 | 9.000000 |
| 25% | 2.263000 | 35.000000 | 12.000000 |
| 50% | 3.461300 | 42.000000 | 12.000000 |
| 75% | 4.970750 | 47.250000 | 14.000000 |
| max | 25.000000 | 60.000000 | 17.000000 |

For women whose husband earn more than the median wage:

```
In [133]: df["wage estimated"] = df["wage estimated"].replace(0, np.nan)
df["wage estimated"] = df["wage estimated"].replace(0, np.nan)
```

The statistics of women salaries whose husband earn more the median are:

| | wage estimated | age women | educ years of schooling |
|-------|----------------|------------|-------------------------|
| count | 214.000000 | 214.000000 | 214.000000 |
| mean | 5.252555 | 41.941121 | 13.135514 |
| std | 2.937960 | 7.338431 | 2.233040 |
| min | 0.303000 | 30.000000 | 7.000000 |
| 25% | 3.576725 | 35.000000 | 12.000000 |
| 50% | 4.543150 | 43.000000 | 12.000000 |
| 75% | 5.863100 | 47.000000 | 15.000000 |
| max | 25.000000 | 58.000000 | 17.000000 |

For women whose husband earn more than the median wage:

```
In [134]: df["wage estimated"] = df["wage estimated"].replace(0, np.nan)
df["wage estimated"] = df["wage estimated"].replace(0, np.nan)
```

The statistics of women salaries whose husband earn more the median are:

| | wage estimated | age women | educ years of schooling |
|-------|----------------|------------|-------------------------|
| count | 214.000000 | 214.000000 | 214.000000 |
| mean | 3.122884 | 42.102804 | 12.182243 |
| std | 3.313609 | 8.108814 | 2.237705 |
| min | 0.128200 | 30.000000 | 9.000000 |
| 25% | 1.697225 | 35.000000 | 12.000000 |
| 50% | 2.890100 | 41.000000 | 12.000000 |
| 75% | 3.313609 | 48.000000 | 15.000000 |
| max | 25.000000 | 60.000000 | 17.000000 |

Question 3

```
In [135]: plt.rcParams["figure.figsize"] = (16, 6)

plt.figure(figsize=(16, 6))
plt.subplot(1, 2, 1)
plt.title('Wage and education')
plt.xlabel('educ years of schooling')
plt.ylabel('wage')

plt.subplot(1, 2, 2)
df['log(wage)'] = np.log(df['wage estimated'])
plt.title('log(wage)')
plt.xlabel('educ years of schooling')
plt.ylabel('log(wage)')
```



On compare les deux variables semblent suivre une loi gaussienne de moyenne approximative 3 et 1.2 respectivement. Ce résultat est normal puisque tout échantillon de donné converge vers une telle loi lorsque n devient assez grand.

Question 4

```
In [136]: r = df[['mothers year of schooling', 'fathers year of schooling', 'mothers year of schooling', 'fathers year of schooling']]
r = r[['mothers year of schooling', 'fathers year of schooling', 'mothers year of schooling', 'fathers year of schooling']]
```

The correlation entre ces deux variables est:

La p-value de rejet

On remarque une corrélation assez forte, ou du moins non-négligeable. Cela peut s'expliquer par le fait que les membres d'un couple sont souvent issus d'un niveau d'études assez similaire à cause de la non-mixité des classes sociales.

Cette corrélation peut-être induite d'un biais de multicollinéarité.

Cependant, puisque la p-value vaut environ 0 (<1%), on rejette l'hypothèse "présence de multicollinéarité" entre les 2 variables.

Question 5

Wage and education / experq / education of the father

```
In [137]: plt.rcParams["figure.figsize"] = (16, 17)

plt.figure(figsize=(16, 17))
plt.subplot(3, 1, 1)
plt.title('Wage and education')
plt.xlabel('educ years of schooling')
plt.ylabel('wage')

plt.subplot(3, 1, 2)
plt.title('Wage and experq')
plt.xlabel('experq')
plt.ylabel('wage')

plt.subplot(3, 1, 3)
plt.title('Wage and fathers year of schooling')
plt.xlabel('fathers year of schooling')
plt.ylabel('wage')
```



Il ne s'agit ici pas d'un effet "toute chose étant égale par ailleurs", puisque pour la variation de chacune des données, les autres données ne sont pas constantes.

Question 6

L'hypothèse est que:

$E(U|X) = 0$

Since: $E(\hat{u}) = E(u) + (x - x')E(u|X)$

$E(\hat{u}) = E(u)$

Le biais de variable censee est le biais lié à la non prise en compte dans le modèle d'une variable importante.

Question 7

```
In [139]: plt.rcParams["figure.figsize"] = (16, 6)

stack = np.column_stack([a,
df['live in SMSA'],
df['educ years of schooling'],
df['experq'],
df['mfefinc'],
df['kids below 6'],
df['kids between 6 and 18']])

mod0 = sm.OLS(df['wage estimated'], stack)
mod0 = mod0.fit()
res0 = mod0.resid
SSR0 = mod0.ssr

plt.hist(res0, label='OLS on wage')

x = np.linspace(-10, 25, 100)
y = norm.pdf(x, 0, 2.758) # we multiply to normalize the 2 curves
plt.plot(x, y, label='gaussian curve')
plt.legend()

plt.show()

print(mod0.summary())
print(SSR0)
```



OLS Regression Results

Dep. Variable: wage estimated R-squared: 0.1405

Model: OLS Adj. R-squared: 0.1314

Method: Least Squares Prob. F-statistic: 1.596-18

Date: Fri, 19 Nov 2021 Log-likelihood: -434.97

No. Observations: 428 AIC: 2194.

DF Model: 421 BIC: 2223.

Covariance type: nonrobust

Omitting 1 variable: nonrobust

coef std err t P>|t| [0.025 0.975]

const -2.2446 0.926 -2.425 0.016 -4.064 -0.425

x1 0.6398 0.071 9.026 0.000 0.498 0.781

x2 0.4617 0.078 5.909 0.000 0.304 0.619

x3 0.0086 0.001 10.258 0.000 0.006 0.010

x4 0.0148 0.015 0.954 0.340 -0.018 0.048

x5 0.0203 0.006 3.061 0.002 0.009 0.031

x6 -0.0789 0.124 -0.632 0.524 -0.315 0.165

Omnibus: 342.215 Durbin-Watson: 2.058

Prob(Omnibus): 0.000 Jarque-Bera (JB): 6326.429

Skew: -0.051 Prob(JB): 2.32e-63

Kurtosis: 28.595 Cond. No. 2.22e+02

Notes: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 2.22e+03. This might indicate that there are strong multicollinearity or other numerical problems.

4886.2143777901156

A part de l'hypothèse, nous pouvons interpréter que les résidus suivent une loi gaussienne centrée en 0.

Cependant, il y a une légère erreur induisant une asymétrie pour des valeurs d'abscisse négatives proches de 0, et de plus, certaines valeurs me semblent incohérentes puisque la courbe de la gaussienne rest pas part: certains points semblent très éloignés des autres - ce qui est inhabituel pour une telle loi.

Question 8

```
In [139]: stack = np.column_stack([a,
df['live in SMSA'],
df['educ years of schooling'],
df['experq'],
df['mfefinc'],
df['kids below 6'],
df['kids between 6 and 18']])

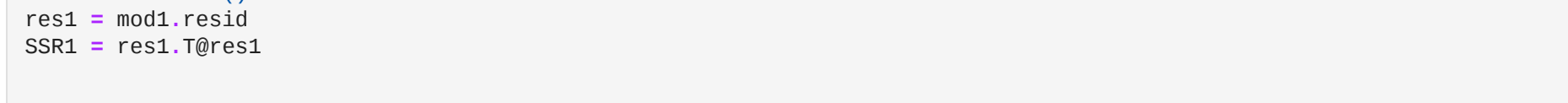
mod1 = sm.OLS(df['log(wage)'], stack)
mod1 = mod1.fit()
res1 = mod1.resid
SSR1 = mod1.ssr

plt.hist(res1, label='OLS on wage')

x = np.linspace(-3, 3, 500)
y = norm.pdf(x, 0, 1.5) # we multiply to normalize the 2 curves
plt.plot(x, y, label='gaussian curve')
plt.legend()

plt.show()

print(mod1.summary())
print(SSR1)
```



OLS Regression Results

Dep. Variable: log(wage) R-squared: 0.1405

Model: OLS Adj. R-squared: 0.1314

Method: Least Squares Prob. F-statistic: 2.37e-12

Date: Fri, 19 Nov 2021 Log-likelihood: -434.97

No. Observations: 428 AIC: 889.1

DF Model: 4 BIC: 911.6

Covariance Type: nonrobust

Omitting 1 variable: nonrobust

coef std err t P>|t| [0.025 0.975]

const -2.2446 0.926 -2.425 0.016 -4.064 -0.425

x1 0.6398 0.071 9.026 0.000 0.498 0.781

x2 0.4617 0.078 5.909 0.000 0.304 0.619

x3 0.0086 0.001 10.258 0.000 0.006 0.010

x4 0.0148 0.015 0.954 0.340 -0.018 0.048

x5 0.0203 0.006 3.061 0.002 0.009 0.031

x6 -0.0789 0.124 -0.632 0.524 -0.315 0.165

Omnibus: 342.215 Durbin-Watson: 2.058

Prob(Omnibus): 0.000 Jarque-Bera (JB): 6326.429

Skew: -0.051 Prob(JB): 2.32e-63

Kurtosis: 28.595 Cond. No. 2.22e+02

Notes: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 2.22e+03. This might indicate that there are strong multicollinearity or other numerical problems.

190.932327389357

On peut donc interpréter que les résidus suivent une loi gaussienne centrée en 0 (ou plutôt juste au dessus de 0), mais cette fois ci les valeurs sont bien mieux distribuées et elles sont cohérentes: il y a beaucoup moins de valeurs extrêmes que dans la question précédente.

Question 9

```
In [132]: SSR1 = res1.TBres1
stack_truncated = np.column_stack([a,
df['live in SMSA'],
df['educ years of schooling'],
df['experq'],
df['mfefinc'],
df['kids below 6'],
df['kids between 6 and 18']])

n, k = np.shape(stack)
mod2 = sm.OLS(df['log(wage)'], stack_truncated)
mod2 = mod2.fit()
res2 = mod2.resid
SSR2 = res2.TBres2

F2 = (SSR2 - SSR1) / SSR1 * (n-k) / k
fisher2 = f.sf(F2, 1, n-k)
print("The p-value is: ", p_value, ", fisher2")

The p-value is: 0.14348799213513744
On ne rejette donc PAS H0 pour aucun des seuls (14-10=4)
```

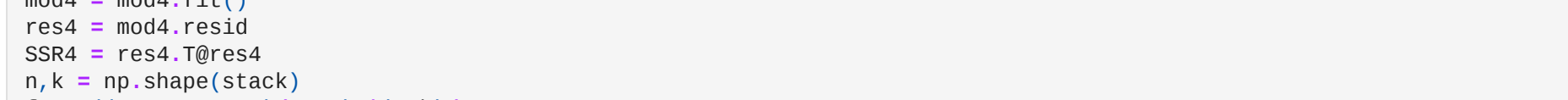
Question 10

```
In [132]: stack3 = np.column_stack([a,
df['live in SMSA'],
df['educ years of schooling'],
df['experq'],
df['mfefinc'],
df['kids below 6'],
df['kids between 6 and 18']])

data3 = df[['log(wage)', '0.81*df['mfefinc']']]
mod3 = sm.OLS(data3, stack3)
mod3 = mod3.fit()
res3 = mod3.resid
SSR3 = res3.TBres3

n, k = np.shape(stack3)
F3 = (SSR3 - SSR1) / SSR1 * (n-k) / k
fisher3 = f.sf(F3, 1, n-k)
print("The p-value is: ", p_value, ", fisher3")

The p-value is: 0.124253902734329
On ne rejette donc pas l'hypothèse
```



OLS Regression Results

Dep. Variable: log(wage) R-squared: 0.1405

Model: OLS Adj. R-squared: 0.1314

Method: Least Squares Prob. F-statistic: 2.37e-12

Date: Fri, 19 Nov 2021 Log-likelihood: -434.97

No. Observations: 428 AIC: 889.1

DF Model: 4 BIC: 911.6

Covariance Type: nonrobust

Omitting 1 variable: nonrobust

coef std err t P>|t| [0.025 0.975]

const -2.2446 0.926 -2.425 0.016 -4.064 -0.425

x1 0.6398 0.071 9.026 0.000 0.498 0.781

x2 0.4617 0.078 5.909 0.000 0.304 0.619

x3 0.0086 0.001 10.258 0.000 0.006 0.010

x4 0.0148 0.015 0.954 0.340 -0.018 0.048

x5 0.0203 0.006 3.061 0.002 0.009 0.031

x6 -0.0789 0.124 -0.632 0.524 -0.315 0.165

Omnibus: 342.215 Durbin-Watson: 2.058

Prob(Omnibus): 0.000 Jarque-Bera (JB): 6326.429

Skew: -0.051 Prob(JB): 2.32e-63

Kurtosis: 28.595 Cond. No. 2.22e+02

Notes: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 2.22e+03. This might indicate that there are strong multicollinearity or other numerical problems.

190.932327389357

On peut donc interpréter que les résidus suivent une loi gaussienne centrée en 0 (ou plutôt juste au dessus de 0), mais cette fois ci les valeurs sont bien mieux distribuées et elles sont cohérentes: il y a beaucoup moins de valeurs extrêmes que dans la question précédente.

Question 11

```
In [138]: stack4 = np.column_stack([a,
df['educ years of schooling'],
df['experq'],
df['mfefinc'],
df['kids below 6'],
df['kids between 6 and 18']])

data4 = df[['log(wage)', '0.81*df['mfefinc']', '0.45*df['live in SMSA']']]
mod4 = sm.OLS(data4, stack4)
mod4 = mod4.fit()
res4 = mod4.resid
SSR4 = res4.TBres4

n, k = np.shape(stack4)
F4 = (SSR4 - SSR1) / SSR1 * (n-k) / k
fisher4 = f.sf(F4, 2, n-k)
print("The p-value is: ", p_value, ", fisher4")

The p-value is: 0.263517405428352
On ne rejette donc pas l'hypothèse à 5%
```

Question 12

```
In [139]: plt.rcParams["figure.figsize"] = (16, 7)

plt.grid()
plt.title('Wage with experience')
both = df.groupby(['educ years of schooling']).mean()
both = both.sort_values(by='educ years of schooling')
plt.plot(both.index, both['wage estimated'])
plt.xlabel('education')
plt.ylabel('wage')
plt.show()
```

