

1. Theory

Question 1

a) $E(\epsilon_i) = \mu_\epsilon \neq 0$

The Formula for $\hat{\beta}_2 = \frac{\sum x_i y_i}{\sum x_i^2}$ according to slide 4 can be rewritten as $\hat{\beta}_2 = \frac{\sum x_i Y_i}{\sum x_i^2}$, since $\sum x_i y_i$ is equal to $\sum x_i Y_i$. For Y_i , one can use the popular regression function and rewrite $\hat{\beta}_2$ as follows:

$$\hat{\beta}_2 = \frac{\sum x_i (\beta_1 + \beta_2 X_i + \epsilon_i)}{\sum x_i^2}$$

This fraction can be separated and the β s can be pulled out:

$$\hat{\beta}_2 = \beta_1 \frac{\sum x_i}{\sum x_i^2} + \beta_2 \frac{\sum x_i X_i}{\sum x_i^2} + \frac{\sum x_i \epsilon_i}{\sum x_i^2}$$

We know that $\sum x_i$ sums up to 0 and that $\sum x_i X_i$ is equal to $\sum x_i^2$, leaving the equation as follows:

$$\hat{\beta}_2 = \beta_2 + \frac{\sum x_i \epsilon_i}{\sum x_i^2}$$

Looking at the expected values, we modify this equation:

$$E[\hat{\beta}_2] = E[\beta_2 + \frac{\sum x_i \epsilon_i}{\sum x_i^2}]$$

$$E[\hat{\beta}_2] = E[\beta_2] + E[\frac{\sum x_i \epsilon_i}{\sum x_i^2}]$$

$$E[\hat{\beta}_2] = \beta_2 + \frac{\sum x_i E[\epsilon_i | X_i]}{\sum x_i^2}$$

$$E[\hat{\beta}_2] = \beta_2 + \frac{\sum x_i \mu_\epsilon}{\sum x_i^2}$$

With $\mu_\epsilon \neq 0$ we can see that $E[\hat{\beta}_2] \neq \beta_2$, meaning that $\hat{\beta}_2$ is biased.

The Formula for $\hat{\beta}_1: \hat{\beta}_1 = \bar{Y} - \hat{\beta}_2 \bar{X}$ includes this biased $\hat{\beta}_2$, meaning $\hat{\beta}_1$ is also biased.

b) $\tilde{X} = 2X$

The formula for $\hat{\beta}$ is, according to slide 5 as follows: $\hat{\beta} = (X'X)^{-1}X'y$

For $\tilde{\beta}$, the formula looks like this:

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

Since we know that $\tilde{X} = 2X$, we can rewrite this formula:

$$\tilde{\beta} = (2X'2X)^{-1}2X'y$$

$$\tilde{\beta} = \frac{1}{4}(X'X)^{-1}2X'y$$

$$\tilde{\beta} = \frac{1}{2}(X'X)^{-1}X'y$$

$$\tilde{\beta} = \frac{1}{2}\hat{\beta}$$

This means that by doubling X , the coefficient $\hat{\beta}$ gets multiplied by $\frac{1}{2}$.

c) $y^* = 2y$

Using the same formula from slide 5 again ($\hat{\beta} = (X'X)^{-1}X'y$), the new equation for β^* looks as follows:

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

Knowing the new y^* being equal to $2y$, we can rewrite this:

$$\beta^* = (X'X)^{-1}X'2y$$

$$\beta^* = 2(X'X)^{-1}X'y$$

$$\beta^* = 2\hat{\beta}$$

This means that by doubling y , the coefficient $\hat{\beta}$ gets multiplied by 2.

d) The results from the previous questions suggest that a linear change in the units in which X and y are measured does not affect the relationship between X and y .

e) From question b) we already know that $\tilde{\beta} = \frac{1}{2}\hat{\beta}$. Therefore, one can rewrite the variance of $\tilde{\beta}$:

$$V(\tilde{\beta}) = V\left(\frac{1}{2}\hat{\beta}\right)$$

Since $\frac{1}{2}$ is a constant multiplied with $\hat{\beta}$, we can pull it out of the variance and square it:

$$V(\tilde{\beta}) = \frac{1}{4}V(\hat{\beta})$$

The formula on slide 5 ($V(\hat{\beta}) = \sigma^2(X'X)^{-1}$) gives the same results:

$$V(\tilde{\beta}) = \sigma^2(\tilde{X}'\tilde{X})^{-1} = \sigma^2(2X'2X)^{-1} = \frac{1}{4}\sigma^2(X'X)^{-1} = \frac{1}{4}V(\hat{\beta})$$

This means that by doubling X , the variance of $\hat{\beta}$ gets multiplied by one quarter.

2. Empirical Application

Question 2

a) Observations

This dataset has 807 observations with 10 variables each.

b) Summary statistics

| Descriptive Statistics | | | | | |
|------------------------|-----|----------|----------|-----|-------|
| Variable | Obs | Mean | Std.Dev. | Min | Max |
| cigs | 807 | 8.686 | 13.722 | 0 | 80 |
| educ | 807 | 12.471 | 3.057 | 6 | 18 |
| age | 807 | 41.238 | 17.027 | 17 | 88 |
| income | 807 | 19304.83 | 9142.958 | 500 | 30000 |
| white | 807 | .879 | .327 | 0 | 1 |
| restaurn | 807 | .247 | .431 | 0 | 1 |

The summary statistics show that all individuals have answered all the questions. The individuals have at least 6 years of schooling. However, the minimum age is 17, meaning that there are probably individuals still being in school, falsifying possible results of the effect of education on cigarettes. All the observed incomes are between 500 and 30'000 with a mean income of 19'304.83. White and restaurn are dummy variables, but the means show that most individuals are white and that only a quarter of the states, in which the individuals live, restrict smoking in restaurants.

c) i) The estimators are computed with the following formulas:

$$\hat{\beta}_2 = \frac{Cov(x,y)}{Va(x)} \text{ and } \hat{\beta}_1 = \bar{Y} - \hat{\beta}_2 \bar{X}$$

The results are: $\hat{\beta}_2 = -.21855212$ and $\hat{\beta}_1 = 11.41203$.

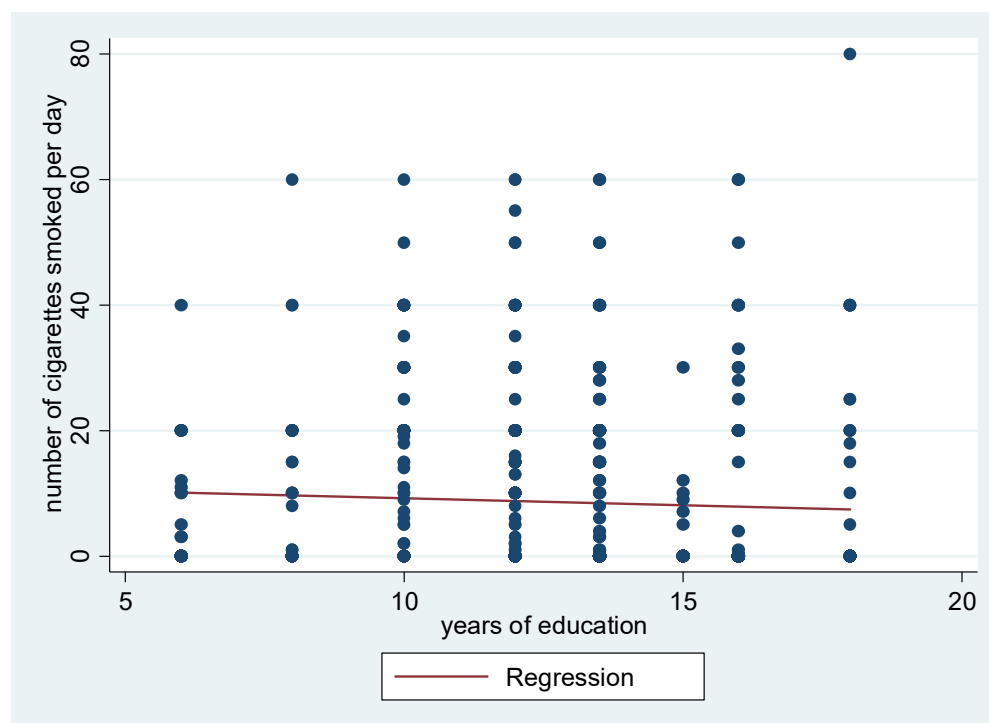
ii) The table below shows the regression results:

| VARIABLES | (1) Reg |
|--------------------------------|---------------------|
| educ | -0.219 (0.158) |
| Constant | 11.41*** (2.029) |
| Observations | 807 |
| R-squared | 0.002 |
| Standard errors in parentheses | |
| *** p<0.01, ** p<0.05, * p<0.1 | |

The estimators in the regression show the exact same results as in 2.c)i).

iii) If the mean-zero error (assumption 2) is satisfied, we know that the estimator is unbiased. The effect of -0.219 means that with one year of additional education, an average individual smokes 0.219 cigarettes fewer per day. On one hand, this effect makes sense, since more education possibly increases awareness of the health risks of smoking. On the other hand, 0.2 cigarettes fewer per day seems like a rather small effect per additional year of education. The effect's standard error is rather large compared to the effect itself, making it insignificant.

iv) $\widehat{cigs}_i = 11.41 - 0.219 \times educ_i$

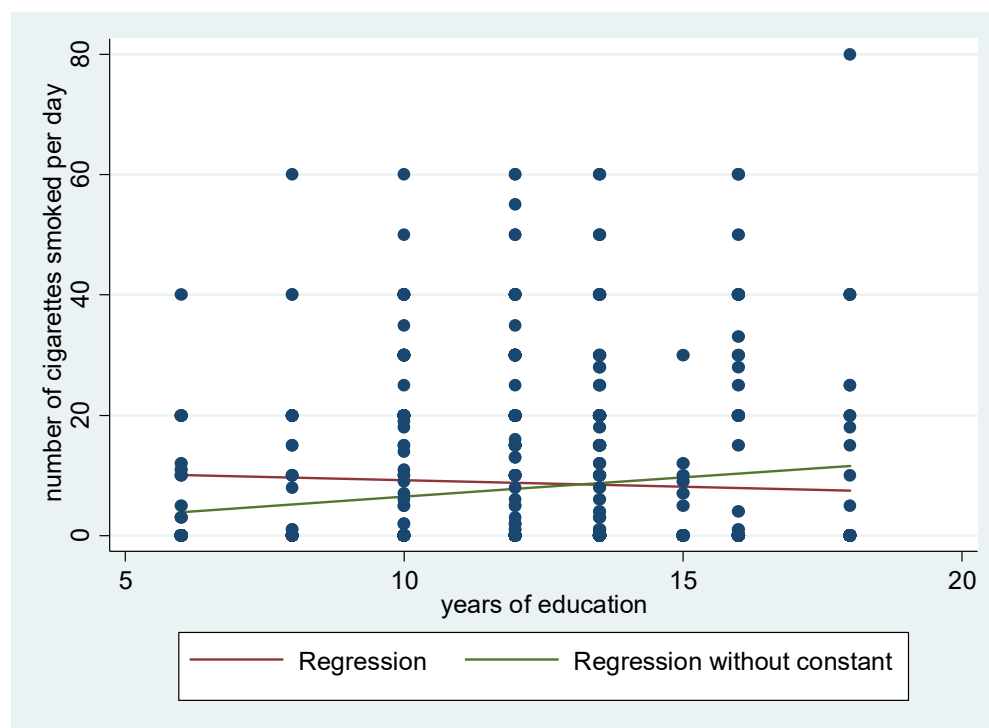


v) The table below shows the regression results:

| VARIABLES | (1) Reg |
|--------------|----------------------|
| educ | 0.645*** (0.0383) |
| Observations | 807 |
| R-squared | 0.260 |

Standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

The effect of education on cigarettes is now positive. Meaning with one additional year of education, an average individual smokes 0.645 cigarettes more. This regression without a constant makes little sense, because it also suggests that an individual with 0 years of education does not smoke at all.



d) i) The table below shows the regression results:

| VARIABLES | (1) Reg |
|--------------|--------------------------|
| educ | -0.452*** (0.162) |
| age | 0.826*** (0.154) |
| age2 | -0.00963*** (0.00168) |
| white | -0.624 (1.456) |
| restaurn | -2.796** (1.104) |
| Constant | 0.669 (3.707) |
| Observations | 807 |
| R-squared | 0.051 |

Standard errors in parentheses

*** p<0.01, ** p<0.05, * p<0.1

The coefficient of white is -0.624. This could lead to the interpretation that on average, individuals of white ethnicity smoke 0.624 cigarettes fewer per day. However, compared to the mean this effect seems rather small and it is also not significantly different from 0. The coefficient of restaurn is -2.796. This leads to the conclusion that individuals who live in a state where smoking in restaurants is restricted smoke on average 2.796 cigarettes fewer per day. This effect is not quite small and the coefficient is significantly different from 0 on a 10%-level.

ii) Marginal effects

In order to find the marginal effect of age, one has to derive the regression with respect to age:

$$\widehat{cigs}' = 0.826 + 2 \times (-0.0096) \times age$$

The calculated marginal effect for 20 years is:

$$\Delta \widehat{cigs}' = 0.826 + 2 \times (-0.0096) \times 20 = 0.442$$

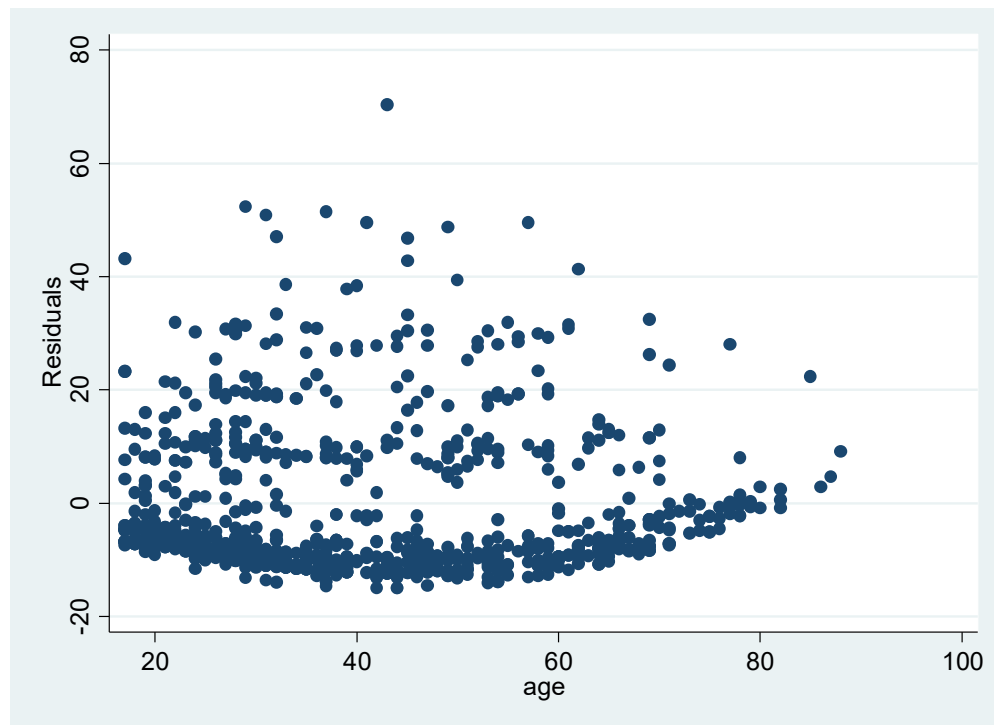
The calculated marginal effect for 40 years is:

$$\Delta \widehat{cigs}' = 0.826 + 2 \times (-0.0096) \times 40 = 0.058$$

The calculated marginal effect for 60 years is:

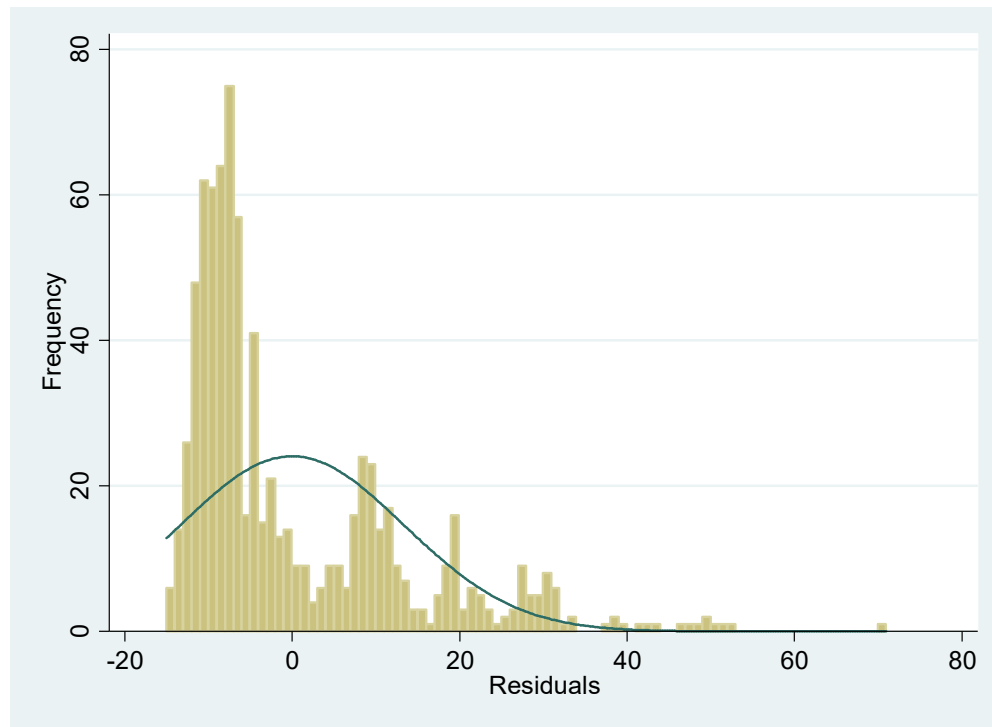
$$\Delta \widehat{cigs}' = 0.826 + 2 \times (-0.0096) \times 60 = -0.326$$

- iii) A) In general, it seems that the variance is constant throughout all ages. However, for older individuals the variance seems to be smaller. This might be explained with relatively fewer observations starting the age of 60. Furthermore, the reason for this model to overestimate the cigarette consumption of individuals older than 60 might be that individuals who consume cigarettes already passed away due to health issues caused by smoking. Therefore, we think that assumption 3 holds.



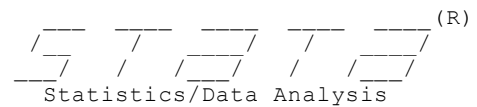
- B) The correlation between the residuals and the lagged residuals is $-.0042332$ and not significantly different from 0. Therefore, we suggest that assumption 4 holds.

C) The distribution of the residuals, compared to a normal distribution, shows a positive skewness, a negative mean and a negative median. This suggests that assumption 5 does not hold.



3. Log-file

See attachment



name: <unnamed>
 log: C:\Users\ramon\Desktop\UZH\Empirical Methods\Problem Sets\Problem Set 1\Stata\log_gm
 log type: smcl
 opened on: 21 Oct 2019, 15:13:29

```
1 .
2 . use "C:\Users\ramon\Desktop\UZH\Empirical Methods\Problem Sets\Problem Set 1\Stata\smoke.dta"
3 .
4 . *8a) How many obs
5 .
6 . display _N
807
```

```
7 .
8 . *8b) Summary statistics for cigs, educ, age, income, white, restaurn
9 .
10 . asdoc sum cigs educ age income white restaurn
    (File Myfile.doc already exists, option append was assumed)
```

| Variable | Obs | Mean | Std. Dev. | Min | Max |
|----------|------------|-----------------|-----------------|------------|--------------|
| cigs | 807 | 8.686493 | 13.72152 | 0 | 80 |
| educ | 807 | 12.47088 | 3.057161 | 6 | 18 |
| age | 807 | 41.23792 | 17.02729 | 17 | 88 |
| income | 807 | 19304.83 | 9142.958 | 500 | 30000 |
| white | 807 | .8785626 | .3268375 | 0 | 1 |
| restaurn | 807 | .2465923 | .4312946 | 0 | 1 |

Click to Open File: [Myfile.doc](#)

```
11 .
12 . *8c)
13 . **i) Compute Beta 1 and 2 (error in the PS, we use 1 and 2, not 0 and 1)
14 . ***B2)
15 . gen COV = 0

16 . correlate educ cigs, covariance
    (obs=807)
```

| | educ | cigs |
|------|-----------------|---------------|
| educ | 9.34624 | |
| cigs | -2.04264 | 188.28 |

```
17 . replace COV = r(cov_12)
    (807 real changes made)
```

```
18 .
19 . egen SD = sd(educ)
```

```
20 . gen VAR = SD^2
```

```
21 .
22 . display COV/VAR
-.21855212
```

```
23 .
24 . ***B1)
```

25 . gen B2 = COV/VAR

26 .

27 . gen mean_cigs = 0

28 . mean(cigs)

Mean estimation Number of obs = **807**

| | Mean | Std. Err. | [95% Conf. Interval] | |
|------|-----------------|-----------------|----------------------|-----------------|
| cigs | 8.686493 | .4830202 | 7.738367 | 9.634619 |

29 . matrix b=e(b)

30 . replace mean_cigs=b[1,1]
(807 real changes made)

31 .

32 . gen mean_educ = 0

33 . mean(educ)

Mean estimation Number of obs = **807**

| | Mean | Std. Err. | [95% Conf. Interval] | |
|------|-----------------|-----------------|----------------------|-----------------|
| educ | 12.47088 | .1076172 | 12.25964 | 12.68212 |

34 . matrix b=e(b)

35 . replace mean_educ=b[1,1]
(807 real changes made)

36 .

37 . gen B1 = mean_cigs - mean_educ*B2

38 .

39 . display B1

11.41203

40 .

41 . **ii) Regression

42 .

43 . reg cigs educ

| Source | SS | df | MS | Number of obs | = | 807 |
|----------|-------------------|------------|-------------------|---------------|---|---------------|
| Model | 359.817074 | 1 | 359.817074 | F(1, 805) | = | 1.91 |
| Residual | 151393.866 | 805 | 188.066914 | Prob > F | = | 0.1670 |
| | | | | R-squared | = | 0.0024 |
| | | | | Adj R-squared | = | 0.0011 |
| Total | 151753.683 | 806 | 188.280003 | Root MSE | = | 13.714 |

| cigs | Coef. | Std. Err. | t | P> t | [95% Conf. Interval] | |
|-------|------------------|-----------------|--------------|--------------|----------------------|-----------------|
| educ | -.2185521 | .1580048 | -1.38 | 0.167 | -.5287022 | .091598 |
| _cons | 11.41203 | 2.028732 | 5.63 | 0.000 | 7.429801 | 15.39426 |

```
44 . outreg2 using "PS1_regression.doc", replace ctitle(Reg)
    PS1_regression.doc
    dir : seeout
```

```
45 .
46 . **iv) Estimates
47 .
48 . graph twoway (lfit cigs educ) (scatter cigs educ)

49 .
50 . **v)
51 .
52 . reg cigs educ, noconstant
```

| Source | SS | df | MS | Number of obs | = | 807 |
|----------|-------------------|------------|-------------------|---------------|---|---------------|
| Model | 55301.1489 | 1 | 55301.1489 | F(1, 806) | = | 283.28 |
| Residual | 157344.851 | 806 | 195.216937 | Prob > F | = | 0.0000 |
| | | | | R-squared | = | 0.2601 |
| | | | | Adj R-squared | = | 0.2591 |
| Total | 212646 | 807 | 263.501859 | Root MSE | = | 13.972 |

| cigs | Coef. | Std. Err. | t | P> t | [95% Conf. Interval] | |
|------|-----------------|-----------------|--------------|--------------|----------------------|-----------------|
| educ | .6447271 | .0383061 | 16.83 | 0.000 | .5695357 | .7199186 |

```
53 . outreg2 using "PS1_regression_noconstant.doc", replace ctitle(Reg)
    PS1_regression_noconstant.doc
    dir : seeout
```

```
54 . twoway (lfit cigs educ) (lfit cigs educ, estopts(noconstant)) (scatter cigs educ)

55 .
56 . *8d)
57 . **i)
58 .
59 . gen age2 = age^2

60 . reg cigs educ age age2 white restaurn
```

| Source | SS | df | MS | Number of obs | = | 807 |
|----------|-------------------|------------|-------------------|---------------|---|---------------|
| Model | 7772.46759 | 5 | 1554.49352 | F(5, 801) | = | 8.65 |
| Residual | 143981.215 | 801 | 179.751829 | Prob > F | = | 0.0000 |
| | | | | R-squared | = | 0.0512 |
| | | | | Adj R-squared | = | 0.0453 |
| Total | 151753.683 | 806 | 188.280003 | Root MSE | = | 13.407 |

| cigs | Coef. | Std. Err. | t | P> t | [95% Conf. Interval] | |
|----------|------------------|-----------------|--------------|--------------|----------------------|------------------|
| educ | -.4515013 | .1615884 | -2.79 | 0.005 | -.768688 | -.1343146 |
| age | .8257641 | .1544737 | 5.35 | 0.000 | .5225431 | 1.128985 |
| age2 | -.009631 | .0016817 | -5.73 | 0.000 | -.012932 | -.00633 |
| white | -.6237386 | 1.45611 | -0.43 | 0.669 | -3.481981 | 2.234504 |
| restaurn | -2.796182 | 1.103552 | -2.53 | 0.011 | -4.962377 | -.6299866 |
| _cons | .6688335 | 3.706849 | 0.18 | 0.857 | -6.607451 | 7.945118 |

```
61 . outreg2 using "PS1_regression2.doc", replace ctitle(Reg)
    PS1_regression2.doc
    dir : seeout
```

```

62 .
63 . **ii)
64 .
65 . reg cigs educ age age2 white restaurn

```

| Source | SS | df | MS | Number of obs | = | 807 |
|----------|-------------------|------------|-------------------|---------------|---|---------------|
| Model | 7772.46759 | 5 | 1554.49352 | F(5, 801) | = | 8.65 |
| Residual | 143981.215 | 801 | 179.751829 | Prob > F | = | 0.0000 |
| | | | | R-squared | = | 0.0512 |
| | | | | Adj R-squared | = | 0.0453 |
| Total | 151753.683 | 806 | 188.280003 | Root MSE | = | 13.407 |

| cigs | Coef. | Std. Err. | t | P> t | [95% Conf. Interval] | |
|----------|------------------|-----------------|--------------|--------------|----------------------|------------------|
| educ | -.4515013 | .1615884 | -2.79 | 0.005 | -.768688 | -.1343146 |
| age | .8257641 | .1544737 | 5.35 | 0.000 | .5225431 | 1.128985 |
| age2 | -.009631 | .0016817 | -5.73 | 0.000 | -.012932 | -.00633 |
| white | -.6237386 | 1.45611 | -0.43 | 0.669 | -3.481981 | 2.234504 |
| restaurn | -2.796182 | 1.103552 | -2.53 | 0.011 | -4.962377 | -.6299866 |
| _cons | .6688335 | 3.706849 | 0.18 | 0.857 | -6.607451 | 7.945118 |

```

66 . mfx, varlist(age age2)

```

Marginal effects after regress
y = Fitted values (predict)
= **8.6864932**

| variable | dy/dx | Std. Err. | z | P> z | [95% C.I.] | | X |
|----------|-----------------|---------------|--------------|--------------|-----------------|-----------------|----------------|
| age | .8257641 | .15447 | 5.35 | 0.000 | .523001 | 1.12853 | 41.2379 |
| age2 | -.009631 | .00168 | -5.73 | 0.000 | -.012927 | -.006335 | 1990.14 |

```

67 .
68 . **iii)
69 . ***A)
70 . reg cigs educ age age2 white restaurn

```

| Source | SS | df | MS | Number of obs | = | 807 |
|----------|-------------------|------------|-------------------|---------------|---|---------------|
| Model | 7772.46759 | 5 | 1554.49352 | F(5, 801) | = | 8.65 |
| Residual | 143981.215 | 801 | 179.751829 | Prob > F | = | 0.0000 |
| | | | | R-squared | = | 0.0512 |
| | | | | Adj R-squared | = | 0.0453 |
| Total | 151753.683 | 806 | 188.280003 | Root MSE | = | 13.407 |

| cigs | Coef. | Std. Err. | t | P> t | [95% Conf. Interval] | |
|----------|------------------|-----------------|--------------|--------------|----------------------|------------------|
| educ | -.4515013 | .1615884 | -2.79 | 0.005 | -.768688 | -.1343146 |
| age | .8257641 | .1544737 | 5.35 | 0.000 | .5225431 | 1.128985 |
| age2 | -.009631 | .0016817 | -5.73 | 0.000 | -.012932 | -.00633 |
| white | -.6237386 | 1.45611 | -0.43 | 0.669 | -3.481981 | 2.234504 |
| restaurn | -2.796182 | 1.103552 | -2.53 | 0.011 | -4.962377 | -.6299866 |
| _cons | .6688335 | 3.706849 | 0.18 | 0.857 | -6.607451 | 7.945118 |

```

71 . rvpplot age

```

```

72 .

```

```

73 . ***B)
74 . predict age_res, residuals
75 . gen age_res1=age_res[_n-1]
    (1 missing value generated)
76 . reg age_res age_res1

```

| Source | SS | df | MS | Number of obs | = | 806 |
|----------|-------------------|------------|-------------------|---------------|---|----------------|
| Model | 2.57931919 | 1 | 2.57931919 | F(1, 804) | = | 0.01 |
| Residual | 143869.779 | 804 | 178.942511 | Prob > F | = | 0.9045 |
| | | | | R-squared | = | 0.0000 |
| | | | | Adj R-squared | = | -0.0012 |
| Total | 143872.358 | 805 | 178.723426 | Root MSE | = | 13.377 |

| age_res | Coef. | Std. Err. | t | P> t | [95% Conf. Interval] | |
|----------|------------------|-----------------|--------------|--------------|----------------------|-----------------|
| age_res1 | -.0042332 | .0352595 | -0.12 | 0.904 | -.0734448 | .0649784 |
| _cons | .0129003 | .4711827 | 0.03 | 0.978 | -.9119931 | .9377938 |

```

77 .
78 . ***C)
79 . hist age_res, frequency normal width(1)
    (bin=86, start=-15.029084, width=1)
80 .
81 .
82 .
    end of do-file

```