

AGEC5213: ECONOMETRIC METHODS
Spring 2019

PROBLEM SET NO. 1
Date Due: February 13, 2019

I. Estimate, Standard Error, t-statistic, and R-squared (5 points)

Researchers are often concerned with how consumers adjust their expenditures on some of addictive commodities (e.g., tobacco, alcohol) as their level of income increases. Consider the following econometric model that relates the proportion a household's income spent on alcohol (ALCOHOL), household's income (INCOME), age of the household head (AGE), and the number of children in the household (NUMKID).

$$ALCOHOL = \beta_1 + \beta_2 \ln(INCOME) + \beta_3 AGE + \beta_4 NUMKID + e$$

A household survey data were used to estimate this model, and the SAS output is kindly presented below. So you don't need to estimate this model by yourself (you can take a nap for a while). One minor problem is that my old laser printer has a problem with its ink cartridge, and as a result some parts of the print-out are left blank.

Dependent variable: ALCOHOL
 Included Observation: 1519

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	0.009052	0.024050	(i)	0.7347
Ln(INCOME)	0.527641	(ii)	6.608620	0.0000
AGE	(iii)	0.00208	-6.962389	0.0000
NUMKID	-0.013282	0.003259	-4.074993	0.0000
R-squared	(iv) $\frac{SSR}{SST}$	Mean dependent var. S.D. dependent var.		0.060596
Adjusted R-squared	0.053047			0.063325
S.E. of Regression (SSR)	(v) $\sqrt{\frac{SSE}{T-K}}$			$t = \frac{b_k - \beta_k}{se(b_k)} = 1.96$
Sum squared error (SSE)	8.752896		$CI = b_k \pm t_{\alpha/2} \cdot SE$	at 95% CI

- (a) Fill in the blank spaces that appear in the output.
- (b) Interpret each of the estimates b_2 , b_3 , and b_4 .
- (c) Compute 95% interval estimates for β_2 and β_3 . What do these interval estimates tell you?
- (d) Test the hypothesis that the household income proportion for alcohol does not depend on the number of children in the household and interpret the test result.

II. Estimates, Standard Errors, and Hypothesis Test (7 points)

Consider the regression model:

$$R^2_{adj} = 1 - \frac{(1-R^2)(T-K)}{(T-K)}$$

$$R^2 = 1 - \frac{(1-R^2)(T-K)}{(T-1)}$$

$$SSR = \sqrt{\frac{SSE}{T-K}}$$

$$R^2 = \frac{SSR}{SST}$$

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

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + e_i; \quad e_i \sim N(0, \sigma^2)$$

$$\hat{\sigma}^2 = \frac{e'e}{N-p} \quad e = y - X\hat{\beta} = \hat{y}$$

Data on this three-variable regression model yields the following results:

$$X'X = \begin{bmatrix} 25 & 0 & 0 \\ 0 & 10 & 20 \\ 0 & 20 & 50 \end{bmatrix}; \quad X'y = \begin{bmatrix} 100 \\ 1 \\ 8 \end{bmatrix}; \quad \sum y_i^2 = 450$$

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

(a) What is the sample size? 25

(b) Estimate parameters of the regression model.

(c) Estimate standard errors of $\hat{\beta}_2$ and $\hat{\beta}_3$.

(d) Test the hypothesis that $\beta_2=0$ against $\beta_2 \neq 0$; and the hypothesis that $\beta_3=1$ against $\beta_3 \neq 1$ (use 5% significance level for both tests).

(e) Compute R^2 and adjusted R^2 (both R^2 should be mean-centered).

$X' \Rightarrow X - \text{transpose}$
find inverse of $X'X$ matrix

III. Elasticity, R-Squared, Adjusted R-Squared, and Hypothesis Test (8 points)

Data on beginning salary (Y) and education in years (X) for 93 employees of Stillwater Bank can be found in the file, HW1-DATA.xls (see our D2L site). Estimate a log-log linear relationship between salary and education, i.e.,

$$\ln Y_i = \beta_0 + \beta_1 \ln X_i + e_i.$$

- (a) Does the coefficient of education have the expected sign? What interpretation can you place on the slope coefficient?
- (b) Predict the starting salary of an individual with 13 years of education.
- (c) What is the interpretation of R-square from your output?
- (d) Compute adjusted R-square and provide justification of using the adjusted R-square.
- (e) Test if the elasticity of salary w.r.t. education is one. Test if all coefficients are zero. List null and alternative hypotheses, the choice of probability distribution, and degree of freedom for these tests.
- (f) In addition to observations on Y and X , the data contains observations on number of months of previous work experience (E) and the number of months that the individual was hired for the Stillwater Bank (T). Run the regression using the same functional form, a double log form. Did b_1 , b_2 , R-squared, adjusted R-squared changed? What can you say about these different regression results?
- (g) Test if elasticities of the number of months of previous work experience (E) and the number of months that the individual was hired for the Stillwater Bank (T) are the same. List null and alternative hypotheses, the choice of probability distribution, and degree of freedom for this test.

March 4

I. Estimate, Standard Error, t-statistics, and R-Squared.

Answers:

(A).

Dependent Variable: Alcohol

Included Observation (N): 1519

Variable	Coefficient (b_k)	Std. Error (SE_{bk})	T-Stat. (t)	Prob. (β_k)
C			-30.17247 (I) $\rightarrow 0.3164$	
Ln(Income)		0.07984133 (II) ✓		
Age	-0.01448177 (III) ✓			
NUMKID				
R-squared	0.53139793 (IV) X			
S.E. of Regression	0.07598472 (V) ✓			

Equations used:

$$t = (b_k - \beta_k) / SE_{bk}$$

S.E. Regression = $\sqrt{\text{Sum Squared Error (SSE)} / (N-k)}$ where, N = number of cases or observation and K is # of variables (4 in our case). (n-k) if degree of freedom.

I. (A).

$$\text{ID} \quad t = \frac{0.009052 - 0.3764}{0.024050} = -30.17247 \quad 0.3764$$

$$\text{II} \quad 6.108620 = \frac{0.527641 - 0}{SE} \quad \text{or } SE = 0.07984133$$

$$\text{III} \quad -6.962389 = \frac{b_k - 0}{0.10208} \quad \text{or } b_k = -0.01448177$$

$$\text{IV} \quad \text{S.E. of Reg} = \sqrt{\frac{\text{SSE}}{T-K}} = \sqrt{\frac{8.752896}{1519-4}} = 0.0759847228$$

$$\text{V} \quad R^2_{\text{Adj}} = 1 - \frac{\text{SSE}/(T-K)}{\text{SST}/(T-1)}$$

$$0.53047 = 1 - \frac{8.752896}{SST} \times \frac{1518}{11815}; SST = 18.67873931$$

$$R^2 = 1 - \frac{\text{SSE}}{\text{SST}} = 1 - \frac{8.752896}{18.67873931} = 0.5313979249$$

$$1 - \frac{R^2_{\text{Adj}} (T-1)}{(T-K)}$$

(B)

Interpretation of b2: with the unit ^{you} percentage increase in the household income keeping other variables constant, the proportion of household income spent on the alcohol increases (+ve sign) by 0.527641 unit. ($\ln \rightarrow \%$)

Interpretation of b3: with the unit increase in the age of person in family ^{on average}, keeping other variables constant, the proportion of household income spent on the alcohol decreases (-ve sign) by 0.0145 unit.

Interpretation of b4: with one additional child in the household and other situation remaining unchanged, the proportion of household income spent on the alcohol decreases (-ve sign) by 0.0133 unit. ^{on average}

C.

Solution c, b2

reg_coeff_b2 = 0.527641

st_err_b2 = 0.00208

Computing 95% CI for b2

ME_b2 = 1.96 * st_err_b2

CI_low_b2 = reg_coeff_b2 - ME_b2

CI_high_b2 = reg_coeff_b2 + ME_b2

CI_b2 = [0.3711519932, 0.6841300068]

$$CI = b_k \pm 1.96 \times SE(b)$$

If repeated sample of same size is drawn from same population, the true value lie between these two values interval for 95% of time.

Solution c, b3

reg_coeff_b3 = -0.01448177

st_err_b3 = 0.00208

Computing 95% CI for b3

ME_b3 = 1.96 * st_err_b3

CI_low_b3 = reg_coeff_b3 - ME_b3

CI_high_b3 = reg_coeff_b3 + ME_b3

CI_b3 = [-0.01855857, -0.01040497]

\Rightarrow Intervals tell that if repeated samples of the same size were taken from the population, the true value of the estimates would lie between the referred intervals 95% of the time.

All the t statistics values that falls inside the range of confidence interval tells that these values are not significantly different than the null value. T-statistics that falls outside this interval are significantly different than null value and thus we can reject null hypothesis at 95% confidence interval. Here both, Ln(Income) and Age are significantly different than their null values as they fall outside the these

confidence interval. So, we reject null hypothesis (here absolute value of t-stat is greater than t-cric) in favor of alternative hypothesis for both b2 and b3 ($\ln(\text{income})$ and Age). The $\ln(\text{age})$ and the income has statistically significant effect in the proportion of household income spent on alcohol.

D. -0.2

Solution d, b4

`reg_coeff_b4 = -0.013282`

`st_err_b4 = 0.003259`

Computing 95% CI for b4

`ME_b4 = 1.96 * st_err_b4`

`CI_low_b4 = reg_coeff_b4 - ME_b4`

`CI_high_b4 = reg_coeff_b4 + ME_b4`

`# CI_b4 = [-0.01966964, -0.00689436]`

$$H_0: \beta = 0$$

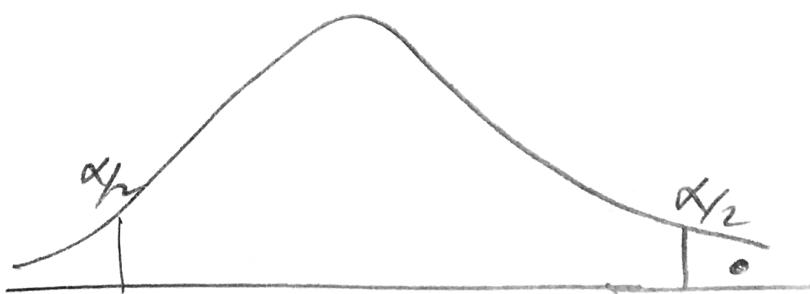
$$H_A: \beta \neq 0$$

$$t = \frac{b_k - \beta}{SE} = \frac{0.527641}{0.07984133} = 6.66.$$

$$t_{\text{cric}} = 1.96 < t_{\text{calc}}$$

The higher and lower critical values at 95% confidence interval is as given above. The calculated t-statistics (-4.074993) does not lies within this range of values in the probability density function and thus the null hypothesis that the household income proportion for alcohol does not depend upon the number of children in the household is rejected. The result suggests that the household income proportion for the alcohol indeed depends upon the number of children in the household.

Ok! In this problem, You should use t-test not interval.



Reject H_0 .

II. Estimates, Standard Errors, and Hypothesis Test

a. What is the sample size? Answer: 25.

b.

B	C	D	E	F	G	H
17						
18	X'X	25	0	0		
19		0	10	20		
20		0	20	50		

Highlight area to put Inverse matrix. $X'X$ Inverse =
 $=\text{MINVERSE}(D18:F20)$. Then hit Ctrl+Shift+Enter at a
 time.

21	(X'X) Inverse	0.04	0	0	100	$\sum Y_t$
22		0	0.5	-0.2	1	
23		0	-0.2	0.1	8	

Determinant of Matrix = mdeterm(array) and press
 enter. Det = MDETERM(D18:F20)

26	Det. Mat X'X	2500				
27						
28	Multiply Matrics:	=MMULT(D22:F24,H22:H24)	$= X'X^{-1}X'y$			
29	Multiply two matrices	4	b1			
30		-1.1	b2			
31		0.6	b3			

B1, b2 and b3 are parameters.

$$(X'X)^{-1} = \begin{bmatrix} 0.04 & 0 & 0 \\ 0 & 0.5 & -0.2 \\ 0 & -0.2 & 0.1 \end{bmatrix} \quad X'y = \begin{bmatrix} 100 \\ 1 \\ 8 \end{bmatrix}$$

$$b = \begin{bmatrix} 4+0+0 \\ 0+0.5-1.6 \\ 0-0.2+0.8 \end{bmatrix} = \begin{bmatrix} 4 \\ -1.1 \\ 0.6 \end{bmatrix} = \begin{matrix} b_1 \\ b_2 \\ b_3 \end{matrix}$$

$$b = (X'X)^{-1}X'y$$

c.

$$SSE = \hat{y}' - \hat{b}'x'y$$

II. (c) standard errors of b_2 & b_3
 $T=25$ $K=3$

$$\underline{\underline{b_2}}$$

$$y'y = \sum_t y_t^2 = 450$$

$$SSE = \hat{e}'\hat{e} = y'y - \hat{b}'x'y$$

$$= 450 - [4 \quad -1.1 \quad 0.6] \begin{bmatrix} 100 \\ 1 \\ 8 \end{bmatrix} = 450 - 400 + 1.1 - 4.8 = 96.3$$

$$SSE = \hat{\sigma}^2 = \frac{\hat{e}'\hat{e}}{T-K} = \frac{SSE}{T-K} = \frac{96.3}{22} = 2.104545$$

$$Var(b) = \hat{\sigma}^2 (X'X)^{-1} = (2.104545) \begin{bmatrix} 0.04 & 0 & 0 \\ 0 & 0.5 & -0.2 \\ 0.1 & -0.2 & 0.1 \end{bmatrix}$$

~~Cov. Cov matrix~~

$$= \begin{bmatrix} 0.0841818 & 0 & 0 \\ 0 & 1.0522725 & -0.420909 \\ 0 & -0.420909 & 0.2104545 \end{bmatrix}$$

$$Se(b_2) = \sqrt{1.0522725} = 1.025803344$$

$$Se(b_3) = \sqrt{0.2104545} = 0.4587532016$$

$$SSE = \hat{e}'\hat{e} - \hat{b}'x'y$$

$$= 450 - [4 \quad -1.1 \quad 0.6] \begin{bmatrix} 100 \\ 1 \\ 8 \end{bmatrix} = 96.3$$

$$\hat{\sigma}^2 = \frac{SSE}{T-K} \quad \& \quad \sigma = \sqrt{\frac{SSE}{T-K}}$$

$$Var(b) = (X'X)\hat{\sigma}^2$$

d.

II d.

 $t = b_k - \beta_k$

$$H_0: \beta_2 = 0$$

$$H_A: \beta_2 \neq 0$$

$$t_{\text{stat}} = \frac{b_k - \beta_k}{SE(b_k)}$$

$$t_{\text{stat}} = \frac{-1.1 - 0}{1.0258} = -1.07$$

You should test using t-value and
t critical value not Interval

$$CI = b_k \pm t_{\alpha/2} \cdot SE$$

$$= -1.1 \pm 1.96 \times 1.0258 = [-3.110568, 0.910568]$$

Since t_{calc} lies within the CI, we failed to
reject null hypothesis

$$H_0: \beta_3 = 1$$

$$H_A: \beta_3 \neq 1$$

$$t_{\text{stat}} = \frac{0.6 - 1}{0.45875} = 0.87193$$

$$CI = b_k \pm t_{\alpha/2} \cdot SE$$

$$= 0.6 \pm 1.967 \times 0.45875 = [-0.80236, 1.50236]$$

Since t_{calc} falls within the confidence
interval, we failed to reject null hypothesis

e.

Q

R²

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} \quad \text{Adj } R^2 = 1 - \frac{\frac{SSE}{(T-K)}}{\frac{SST}{(T-1)}}$$

$$X'y = \begin{bmatrix} \sum y_t \\ \sum x_{1t} \cdot y_t \\ \sum x_{2t} \cdot y_t \end{bmatrix} = \begin{bmatrix} 100 \\ 1 \\ 8 \end{bmatrix}$$

\hat{y} or \bar{y} Average of y_t

$$\sum y_t = 100, \quad \sum y_t^2 = 450 \quad \hat{y} = \frac{\sum y_t}{n} = \frac{100}{25} = 4$$

$$\begin{aligned} SST &= \sum_n (y_t - \bar{y})^2 \\ &= \sum_n (y_t^2 - 2y_t\bar{y} + \bar{y}^2) = \sum_n (y_t^2 - 2y_t \cdot 4 + 4^2) \\ &= \sum_n y_t^2 - 8\sum y_t + 16\sum_n \\ &= 450 - 8 \cdot 100 + 16 \cdot 25 \\ &= 50 \end{aligned}$$

$$SSE = 46.3$$

$$SSR = SST - SSE = 50 - 46.3 = 3.7$$

$$R^2 = \frac{SSR}{SST} = \frac{3.7}{50} = 0.074$$

$$R_{\text{Adj}}^2 = 1 - \frac{\frac{SSE}{(T-K)}}{\frac{SST}{(T-1)}} = 1 - \frac{\frac{(46.3)/22}{50/24}}{\frac{50/24}{50/24}} = 1 - 0.01018$$

✓

$$= -0.010182$$

III. Elasticity, R-Squared, Adjusted R-Squared, and Hypothesis Test

a.

The SAS System

The REG Procedure
Model: MODEL1
Dependent Variable: Iny

Number of Observations Read	94
Number of Observations Used	93
Number of Observations with Missing Values	1

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	0.22901	0.22901	15.94	0.0001
Error	91	1.30703	0.01436		
Corrected Total	92	1.53604			

Root MSE	0.11985	R-Square	0.1491
Dependent Mean	8.58961	Adj R-Sq	0.1397
Coeff Var	1.39524		

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	7.95996	0.15817	50.32	<.0001
InX	1	0.25108	0.06288	3.99	0.0001

a. I expect the rise in salary with the higher level of education i.e. positive relationship between these two variables. So, yes, the education has the expected sign with the salary. With the unit percentage change (say, increase) in the education level, the change in the salary (increase in this case) is by 0.25108% 

b. $\ln(y) = 7.95996 + 0.25108 \ln(x)$

When $x = 13$, $\ln(y) = 7.95996 + 0.25108 \ln(13) = 8.603967$.

$$\hat{Y} = \text{EXP}(8.603967) = 5453.252$$

c. The value of $R^2 = 0.1491$ (14.91%) means the $\ln(x)$ (natural log of education in years) explains 14.91% of variation in the $\ln(y)$ (natural log of salary).

d. -0.5

$$R^2 = 1 - \frac{\text{SSE}/(T-k)}{\text{SST}/(T-1)}$$

$$\text{SSE} = 1 - \frac{\text{SST}}{\text{SSE}} \times \frac{(T-1)}{(T-k)}$$

$$= 1 - \frac{1.30703}{1.53604} \times \frac{92}{91}$$

$$= 1 - 0.860259477$$

$$= 0.13974$$

R^2 squared explains the variation to all the independent variables in the model. So, as we add more variables into the model, we will get higher R^2 -squared value. However, adjusted R^2 -squared only explains the variation in the model due to significant variables in the model. So, it is not possible to increase the value of adjusted R^2 -squared just by adding new variable into the model.

\bar{R}^2 value is corrected for the number of regressors (independent variable) included. (You should mention the relationships between the number of independent variable and value of R^2 and \bar{R}^2)

The value of \bar{R}^2

e.

-0.5

The SAS System

The REG Procedure Model: MODEL1

Test 1 Results for Dependent Variable Iny				
Source	DF	Mean Square	F Value	Pr > F
Numerator	1	2.03761	141.87	<.0001
Denominator	91	0.01436		

Null Hypothesis: $H_0: \beta_2 = 1$

Alternative Hypothesis $H_a: \beta_2 \neq 1$

Choice of Probability Distribution: F Distribution

Degree of Freedom = f (1, 91)

The result shows that the elasticity of salary w.r.t. education is not one (Pr < 0.001), reject null hypothesis.

The SAS System

The REG Procedure Model: MODEL1

Test 2 Results for Dependent Variable Iny				
Source	DF	Mean Square	F Value	Pr > F
Numerator	1	36.37462	2532.53	<.0001
Denominator	91	0.01436		

Null Hypothesis: $H_0: \beta_1 = 0$

Alternative Hypothesis $H_a: \beta_1 \neq 0$

Choice of Probability Distribution: F Distribution

Degree of Freedom = f (1, 91)

We reject null hypothesis that intercept is zero ($P < 0.0001$).

The SAS System

The REG Procedure
Model: MODEL1

Test 3 Results for Dependent Variable lny				
Source	DF	Mean Square	F Value	Pr > F
Numerator	1	0.22901	15.94	0.0001
Denominator	91	0.01436		

Null Hypothesis: $H_0: \beta_2 = 0$;Alternative Hypothesis $H_a: \beta_2 \neq 0$

$$H_0: \beta_1 = 0, \beta_2 = 0$$

$$H_1: \beta_1 \neq 0 \text{ and/or } \beta_2 \neq 0$$

Choice of Probability Distribution: F Distribution

Degree of Freedom = f (1, 91)

The result shows that the lnx is not zero ($Pr < 0.001$). We can reject null hypothesis.

f.

The SAS System**The REG Procedure**

Model: MODEL1

Dependent Variable: Inv

Number of Observations Read	94
Number of Observations Used	90
Number of Observations with Missing Values	4

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	3	0.44972	0.14991	14.43	<.0001
Error	86	0.89354	0.01039		
Corrected Total	89	1.34326			

Root MSE	0.10193	R-Square	0.3348
Dependent Mean	8.59773	Adj R-Sq	0.3116
Coeff Var	1.18556		

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	7.72596	0.14504	53.27	<.0001
InX	1	0.25183	0.05350	4.71	<.0001
Ine	1	0.03694	0.00990	3.73	0.0003
Int	1	0.03354	0.01177	2.85	0.0055

Yes, b_1 , b_2 , R-squared and Adjusted R-squared value changed. The addition and removal of independent variables in the analysis changes the model that we are developing or using to predict the dependent variable which also changes these values.

The addition of two more variables in the model increased the value of adjusted R-squared which signifies that addition of these variables help to develop better model than before.

g.

The SAS System

The REG Procedure Model: MODEL1

Test 1 Results for Dependent Variable Inv				
Source	DF	Mean Square	F Value	Pr > F
Numerator	1	0.00050707	0.05	0.8257
Denominator	86	0.01039		

Null Hypothesis: $H_0: \beta_2 - \beta_3 = 0$; $H_1: \beta_3 = \beta_4$

Alternative Hypothesis $H_a: \beta_2 - \beta_3 \neq 0$ $H_a: \beta_3 - \beta_4 \neq 0$

Choice of Probability Distribution: F Distribution

Degree of Freedom = f (1, 86)

We failed to reject null hypothesis based on the P value (0.8257) which is greater than 0.05 at 95% confidence interval.

$$F_{1,86} = 3.94$$

need conclusion,,

SAS Code:

Import Data and Open Project:

PROC IMPORT OUT= WORK.bm

DATAFILE= "C:\Users\casnrlab_agh128\Desktop\EconHW\HW1-DATA.xls"

DBMS= EXCEL REPLACE;

GETNAMES=YES;

DATAROW=2;

RUN;

dbms = excel replace;

range = sheet1\$

getnames = yes;

mixed = no;

scantext = yes;

scantime = yes;

run;

data bm; set bm;

lny = log (y);

lnx = log(x);

run;

a. # e.

proc reg data = bm;

model lny = lnx;

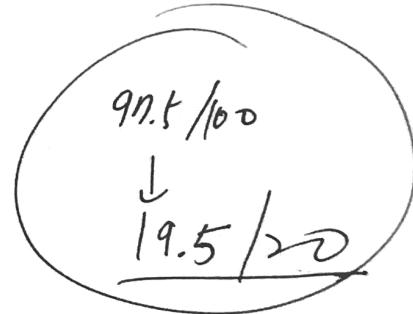
test lnx = 1;

test intercept = 0;

test lnx = 0;

run;

proc print;



```
run;
```

```
#f #g
```

```
data bm; set bm;
```

```
lny = log(y);
```

```
lnx = log(x);
```

```
lne = log(e);
```

```
lnt = log(t);
```

```
proc reg data = bm;
```

```
model lny = lnx lne lnt;
```

```
test lne-lnt = 0;
```

```
run;
```

```
proc print;
```

```
run;
```