BIOSTAT 200B HW5

Part A.

1. For the regression of logdocper on logpopdens bedp1000 hsgrad poverty unemp pcinck, provide an interpretation of each of the regression coefficients.

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Variance Inflation
Intercept	1	-2.98003	0.31428	-9.48	<.0001	0
logpopdens	1	0.04832	0.01573	3.07	0.0023	1.56293
bedp1000	1	0.15378	0.00854	18.01	<.0001	1.36750
hsgrad	1	0.01711	0.00332	5.16	<.0001	2.53581
poverty	1	0.04025	0.00518	7.76	<.0001	2.73125
unemp	1	-0.02622	0.00815	-3.22	0.0014	1.70162
pcinck	1	0.06515	0.00556	11.72	<.0001	2.38561

The regression coefficients of the regression model are shown above.

The coefficient of **logpopdens** is 0.04832, which means when the natural logarithm of population density increases one unit, the natural logarithm of doctors per capita times 1000 will increase by 0.04832. That is to say, for any 1% increase in the population density increase, the expected ratio of doctors per capita times 1000 will be $(1.01)^{0.04832}$ = 1.00048 (0.048% increase in doctors per capita).

The coefficient of **bedp1000** is 0.15378, which means when hospital beds per capita times 1000 increases one unit, the natural logarithm of doctors per capita times 1000 will increase by 0.15378. That is to say, for a 1000 increase in hospital beds per capita, we expect to see about 16.62% increase in doctors per capita times 1000, since $e^{0.15378} = 1.1662$.

The coefficient of **hsgrad** is 0.01711, which means when the percent of adult population (25 years and older) who completed 12 or more years of school increases one unit, the natural logarithm of doctors per capita times 1000 will increase by 0.01711. That is to say, for one percent increase in adult population (25 years and older) who completed 12 or more years of school, we expect to see about 1.73% increase in doctors per capita times 1000, since $e^{0.01711} = 1.0173$.

The coefficient of **poverty** is 0.04025, which means when the percent of 1990 population with income below poverty level increases one unit, the natural logarithm of doctors per capita times 1000 will increase by 0.04025. That is to say, for one percent increase in population with income below poverty level, we expect to see about 4.11% increase in doctors per capita times 1000, since $e^{0.04025} = 1.0411$.

The coefficient of **unemp** is -0.02622, which means when the percent of 1990 labor force that was unemployed increases one unit, the natural logarithm of doctors per capita times

1000 will decrease by 0.02622. That is to say, for one percent increase in labor force that was unemployed, we expect to see about 2.59% decrease in doctors per capita times 1000, since $e^{-0.02622} = 0.9741$.

The coefficient of **pcinck** is 0.06515, which means when the per capita income divided by 1000 increases one unit, the natural logarithm of doctors per capita times 1000 will increase by 0.06515. That is to say, for a 1000 increase in the per capita income, we expect to see about 6.73% increase in doctors per capita times 1000, since $e^{0.06515}$ = 1.0673.

The intercept is -2.98003, which means when all the other predictor variables is set to 0, the mean natural logarithm of doctors per capita times 1000 is -2.98003. (doctors per capita times 1000 is $e^{-2.98003} = 0.0508$)

2. For the predictor pcinck, verify that the SE of the regression coef is equal to the formula involving R^2_j given in lecture, ie, find each of the quantities in the SE and calculate the SE, verifying it is equal to the SE given in the SAS output for the regression.

The SE of the regression coefficient for the predictor **pcinck** is 0.00556. According to the lecture,

$$Var(\hat{\beta}_j) = \frac{1}{(1 - R_j^2)} \frac{\sigma^2}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2} = \frac{1}{(1 - R_j^2)} \frac{\sigma^2}{(n-1)s_j^2}$$

$$Var(\widehat{\beta}_{j}) = \frac{1}{(1 - 0.5808)} * \frac{0.0937}{(440 - 1) * 4.059192^{2}}$$

$$Var(\widehat{\beta}_{j}) = SE^{2} = 0.00003090114$$

$$SE = 0.00555887938 \approx 0.00556$$

proc reg data=cdi;

model pcinck = logpopdens bedp1000 hsgrad

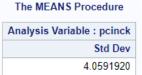
poverty unemp;

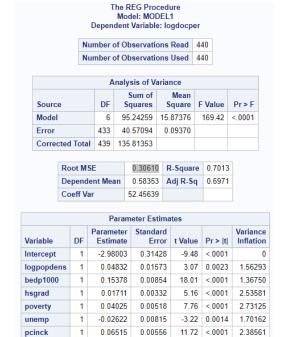
run;

proc means data=cdi std;

var pcinck;

run;

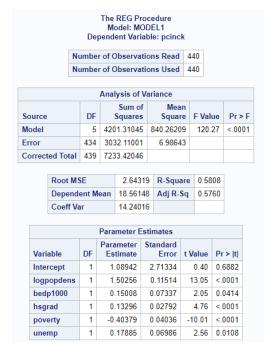




We can get R_i^2 by regress other predictor variables on the predictor **pcinck**, which is 0.5808 and shown in the ANOVA table beside.

Form the ANOVA table of the original regression model, we get σ^2 equals to MSE (0.0937).

To acquire the standard deviation of **pcinck**, we can do proc means in SAS and get the result 4.0591920.



3. In general, would you expect VIFs to increase or decrease as more predictors are added to a model? Provide a justification for your answer. Illustrate by providing some example regressions using the CDI dataset.

I expect VIF to increase as more predictors are added to a model. Since $VIF_j = \frac{1}{1 - R_i^2}$, we can expect VIF_j increase while R_j^2 is high, which is high only when the other predictors can explain most of the variance in the regression on variable X_j . The more predictors are added to a model, the more variance in the regression on variable X_j can be explained by those predictors. That is to say, when more predictors are added to a model, R_j^2 will be higher due to the higher possibility that those predictors are correlated and can explain each other, and VIF_j will increase consequently.

proc reg data=cdi; model logdocper = logpopdens /vif; run;

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Variance Inflation
Intercept	1	-0.77646	0.12220	-6.35	<.0001	0
logpopdens	1	0.22824	0.02013	11.34	<.0001	1.00000

proc reg data=cdi;

model logdocper = logpopdens bedp1000 pcinck /vif;

run;

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Variance Inflation
Intercept	1	-1.45360	0.09244	-15.73	<.0001	0
logpopdens	1	0.05622	0.01669	3.37	0.0008	1.48601
bedp1000	1	0.17589	0.00834	21.08	<.0001	1.10297
pcinck	1	0.05712	0.00465	12.28	<.0001	1.41043

proc reg data=cdi;

model logdocper = logpopdens bedp1000 hsgrad bagrad poverty unemp pcinck crmp1000 pop pop18 pop65 totalinc /vif;

run;

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Variance Inflation
Intercept	1	-2.26170	0.41397	-5.46	<.0001	0
logpopdens	1	0.00950	0.01602	0.59	0.5535	1.99481
bedp1000	1	0.15038	0.00860	17.48	<.0001	1.70840
hsgrad	1	0.00157	0.00378	0.41	0.6787	4.05862
bagrad	1	0.02540	0.00431	5.90	<.0001	6.26396
poverty	1	0.01706	0.00607	2.81	0.0052	4.61291
unemp	1	-0.01119	0.00787	-1.42	0.1559	1.95175
pcinck	1	0.04033	0.00885	4.56	<.0001	7.43660
crmp1000	1	0.00147	0.00070906	2.08	0.0384	1.79836
рор	1	3.250915E-7	2.004061E-7	1.62	0.1055	83.94670
pop18	1	0.01341	0.00511	2.62	0.0090	2.64898
pop65	1	0.02095	0.00463	4.52	<.0001	1.97266
totalinc	1	-0.00001148	0.00000961	-1.19	0.2328	88.42120

As the example we shown above, the VIF for the first model is 1 because there is only one predictor variable. When there are three predictor variables in the second model, we get those VIFs to be less than 1.5, since those variables are not highly correlated; thus, R² is not high and so is VIF. We put a lot of predictors variables in the third model, where we can notice that some corresponding VIFs are really high, which shows collinearity in this regression model. This is because we add too many correlated variables in this model, and those variables could be considered to be removed because they can be explained by other predictor variables.

```
Biosener 200B HWS
 Pare B. Cease squares using macrix algebra
  1. I-H= I-X(X'X)-'X'
       (I-H)'= (I-X(X'X)-'X')'
                = I'- (x')'((x'x)-1)'X'
                = I - X(x'x) x' = I - H = I - H is symmetric
      (I-H)(I-H)= I2- IH-HI+H2
                    = I - 2H+H2
                    = I - 2H + (X(X'X) - X') [X (X'X) - X']
                    = I - 2H+ X(X'X) - X'X X - (X') - X'
                   = I - 2H + X(XX) 1X'
                   = I-VH+H = I-H = I-H is idempotent
2 I- tJ = I- [ + ]
     (I-hJ)'= I'-hJ'
                = I- nJ, : I is symmetric and J= [ ] is also symmetric
               => I-hJ is symmetric
    (エーカー)(エーカー) = エューカーナーカーナーカー
                       = I - \frac{2}{n}J + \frac{1}{n^2}J^2
= I - \frac{2}{n}J + \frac{1}{n^2}(\frac{1}{n}J^2)
= I - \frac{2}{n}J + \frac{1}{n^2}(\frac{1}{n}J^2)
                       = I - = J + - J = I - nJ = I - nJ is idempotent
3. (a) E(Y) = E(X\beta + \varepsilon) = E(X\beta) + E(\varepsilon)
                                                       (XB is a vector of constants)
                                = X\beta + 0 = X\beta
      Var(Y) = Var(X\beta + E) = Var(E) = 6^2I
 (b) The normal equation for least squares estimation of regression parameters:
              \begin{cases} n \hat{\beta}_0 + \hat{\beta}_1 \sum X_i = \sum Y_i \\ \hat{\beta}_0 \sum X_i + \hat{\beta}_1 \sum X_i^2 = \sum X_i Y_i \end{cases}
```

= HY

$$E(\hat{\gamma}) = E(HY) = E(X(XX)XX\beta)$$

$$= E(X\beta) = X\beta \quad (: X\beta : a vector of constants)$$

$$Var(\hat{\gamma}) = Var(HY)$$

$$= H62IH' \quad (: H' = H : B : HH = H : HH' = H)$$

$$= 8^{2}H$$

$$4. \quad yi = 60 + \beta i Xi + \epsilon i$$

$$(a) \quad x_{2nx} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}^{n}$$

$$(b) \quad x_{2nx} = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}^{n}$$

$$(b) \quad x_{2nx} = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}^{n}$$

$$(c) \quad x_{2nx} = \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ -1 & n \end{bmatrix}^{n}$$

$$(c) \quad x_{2nx} = \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n}$$

$$(c) \quad x_{2nx} = \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n}$$

$$(c) \quad x_{2nx} = \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n}$$

$$(c) \quad x_{2nx} = \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n}$$

$$(c) \quad x_{2nx} = \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n}$$

$$(c) \quad x_{2nx} = \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n}$$

$$(c) \quad x_{2nx} = \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n}$$

$$(c) \quad x_{2nx} = \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}^{n} \begin{bmatrix} 1 & -1 \\ 0$$

(c) $\beta = (x'x)^{-1}x'y = \begin{pmatrix} \frac{1}{h} - \frac{1}{h} \\ -\frac{1}{h} & \frac{2}{h} \end{pmatrix} \begin{pmatrix} \frac{2}{h} u_1 + \frac{2}{h} v_1 \\ \frac{2}{h} v_2 \end{pmatrix} = \begin{pmatrix} \frac{1}{h} \frac{2}{h} u_1 + \frac{1}{h} \frac{2}{h} v_1 - \frac{1}{h} \frac{2}{h} v_1 - \frac{1}{h} \frac{2}{h} v_1 \\ -\frac{1}{h} \frac{2}{h} u_1 - \frac{1}{h} \frac{2}{h} v_1 + \frac{1}{h} \frac{2}{h} v_1 \end{pmatrix} = \begin{pmatrix} \overline{u} \\ \overline{v} - \overline{u} \end{pmatrix}$ Bo means when Xi is set to be O (i.e. parlings are females), the mean if outcome varlakes is ti; when Xi is set to be 1 (i.e. patients are males), the expected mean of outcome variables is v-in, which is shown as for here. (d) $H = \chi(\chi \chi)^{-1} \chi' = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{h} & -\frac{1}{h} \\ -\frac{1}{h} & \frac{2}{h} \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{h} & -\frac{1}{h} \\ 0 & \frac{1}{h} \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}$

$$H = \begin{bmatrix} \frac{1}{n} - \frac{1}{n} & 0 - \cdots & 0 \\ \frac{1}{n} - \frac{1}{n} & 0 - \cdots & 0 \\ \frac{1}{n} - \frac{1}{n} & 0 - \cdots & 0 \\ 0 - \cdots & 0 & \frac{1}{n} - \frac{1}{n} \end{bmatrix}_{n \times n}^{n}$$

$$\hat{Y} = HY = \begin{cases} \frac{1}{n_{nxn}} & O_{nxn} \\ O_{nxn} & I_{nxn} \\ O_{nxn} & I_{nxn} \end{cases} \begin{cases} u_1 \\ u_2 \\ u_3 \\ v_4 \end{cases} = \begin{cases} \frac{1}{n_{xx}} u_2 \\ \frac{1}{n_{xx}} u_3 \\ \frac{1}{n_{xx}} u_4 \\ \frac{1}{n_{xx}} u_4 \end{cases} = \begin{cases} \frac{1}{n_{xx}} u_4 \\ \frac{1}{n_{xx}} u_4 \\ \frac{1}{n_{xx}} u_4 \\ \frac{1}{n_{xx}} u_4 \\ \frac{1}{n_{xx}} u_4 \end{cases} = \begin{cases} \frac{1}{n_{xx}} u_4 \\ \frac{1}{n_$$

$$|f| \leq E = e'e = (Y - \hat{Y})'(Y - \hat{Y})$$

$$= \left[u_1 - \bar{u} \dots u_{n-\bar{u}} \cdot \bar{v} - \bar{v} \cdot \bar{v} - \bar{v}\right] \left\{\begin{array}{l} u_1 - \bar{u} \\ v_1 - \bar{v} \\ v_{n-\bar{v}} \end{array}\right\} = \sum_{\bar{v} = 1}^{n} (u_{\bar{v}} - \bar{u}) + \sum_{\bar{v} = 1}^{n} (v_{\bar{v}} - \bar{v})$$

$$MSE = S^{2} = \hat{G}^{2} = \frac{SSE}{2n-2}$$

$$= \frac{\sum_{i=1}^{N} (u_{i} - \bar{u}) + \sum_{i=1}^{N} (V_{i} - \bar{v})}{2n-2}$$

$$= \left(\frac{\sum_{i=1}^{N} (u_{i} - \bar{u}) + \sum_{i=1}^{N} (V_{i} - \bar{v})}{n-1}\right) \times \frac{1}{2}$$

$$= \frac{1}{2} \left(S_{u}^{2} + S_{v}^{2}\right)$$

(g) Ho. Mfensle = Mmale, assume that
$$6^2$$
 female = 6^2 male $6^2 = \frac{(n-1)5u^2 + (n-1)5v^2}{(n-1)+(n-1)} = \frac{1}{2}(5u^2 + 5v^2)$

(h) The estimates for variance in (f) and (g) are the same, = (50 + 52)

This is because we have the same sample size for each groups in this case.