

Question 1

The following model can be used to study whether campaign expenditures affect election outcomes:

$$\text{voteA} = \beta_0 + \beta_1 \ln[\text{expendA}] + \beta_2 \ln[\text{expendB}] + \beta_3 \text{prystrA} + u$$

$$\text{voteA} = 45.1 + 6.08 \ln[\text{expendA}] - 6.62 \ln[\text{expendB}] + 0.152 \text{prystrA} + u \quad n=173$$

1. What is the interpretation of β_1 ?

Soln: When the campaign expenditure of candidate A increases/decreases by 1% then there is $\beta_1\%$ increase/decrease in receiving vote, i.e., when there is 1% increase in expendA then there is 6.08% increase in voteA received.

2. In terms of the parameters, state the null hypothesis that a 1% increase in A's expenditures is offset by a 1% increase in B's expenditures.

Soln: Null hypothesis: $H_0: \beta_1 = -\beta_2 \rightarrow$ if there is 1% increase in expenditure A and 1% increase in expenditure B then percentage change in votes received is not changed

Alternate hypothesis: $H_1: \beta_1 \neq -\beta_2 \rightarrow$ if there is 1% increase in expenditure A and 1% increase in expenditure B then percentage change in votes received is changed

3. Estimate the given model using the data in the vote1 table and report the results in usual form. Do A's expenditures affect the outcome? What about B's expenditures? Can you use these results to test the hypothesis in part 2?

Soln:

R-Code:

```
modlea <- lm(voteA~log(expendA)+log(expendB)+prtystrA,data=vote)
summary(modlea)
tidyw(modlea)
```

$$\text{voteA} = 45.1 + 6.08 \ln[\text{expendA}] - 6.62 \ln[\text{expendB}] + 0.152 \text{prystrA} + u \quad n=173, R\text{-squared}=0.7925$$

The coefficient of $\ln[\text{expendA}]$ and $\ln[\text{expendB}]$ are significant at 0.1% level and the t-statistics for expendA and expendB are 15.915 and -17.461 which is high. Since the coefficients of $\ln(\text{expendA})$ and $\ln(\text{expendB})$ are opposite sign we can test the hypothesis in part 2

4. Estimate a model that directly gives the t-statistic for testing the hypothesis in part 2

Soln:

R-code:

```
summary(lm(voteA~log(expendA)+I(log(expendB)-log(expendA))+prtystrA,data=vote))
```

$$\theta_1 = \beta_1 + \beta_2 \rightarrow \beta_1 = \theta_1 - \beta_2$$

$$\text{voteA} = \beta_0 + (\theta_1 - \beta_2) \ln[\text{expendA}] + \beta_2 \ln[\text{expendB}] + \beta_3 \text{prystrA} + u$$

$$\text{voteA} = \beta_0 + \theta_1 \ln[\text{expendA}] + \beta_2 (\ln[\text{expendB}] - \ln[\text{expendA}]) + \beta_3 \text{prystrA} + u$$

$$\text{voteA} = 45.1 - 0.534 \ln[\text{expendA}] - 6.62 \beta_2 (\ln[\text{expendB}] - \ln[\text{expendA}]) + 0.152 \text{prystrA} + u$$

5. What do you conclude? (Use a two-sided alternative.)

Soln: Since $\ln(\text{expendA})$ is not significant any all the levels we cannot reject the null hypothesis in part 2

Question 2

1. Using the model: $\ln[\text{salary}] = \beta_0 + \beta_1 \text{LSAT} + \beta_2 \text{GPA} + \beta_3 \ln[\text{libvol}] + \beta_4 \ln[\text{cost}] + \beta_5 \text{rank} + u$, state and test the null hypothesis that the rank of law schools has no ceteris paribus effect on median starting salary

Soln:

R-code:

```
model <- lm(log(salary)~ LSAT + GPA +log(libvol) +log(cost)+rank ,data=lawsch)
tidy(model)
```

H0: $\beta_5=0$

```
log(salary)=8.34+ 0.00470 LSAT + 0.248 GPA +0.0950 log(libvol) +0.0376 log(cost) + -0.00332rank,
data=lawsch)
```

Since the coefficient of rank are significant at 0.1% level ($p=1.12e-16$), we can reject the null hypothesis

2.Are features of the incoming class of students—namely, LSAT and GPA—individually or jointly significant for explaining salary? (Be sure to account for missing data on LSAT and GPA.)

Soln:

R Code:

```
lawsch2<-na.omit(lawsch) #removing the na values
modela <- lm(log(salary)~ LSAT+GPA+log(libvol)+log(cost)+rank,data=lawsch2)
tidy(modela)
summary(modela)
modelb <- lm(log(salary)~ log(libvol)+log(cost)+rank,data=lawsch2)
summary(modelb)
anova(modela,modelb)
```

I have omitted the na values in dataset

```
ln(salary) = 7.85+ 0.00683 LSAT+0.233 GPA+0.106 log(libvol)+ 0.0494 log(cost)- 0.00291 rank,
data=lawsch2)
```

In the above model we see the LSAT has t-statistics = 1.23 and P value = 0.223 which is not significant and GPA with t-statistics=2.025 and p-value=0.0460 which is significant at 10%. The joint significance has the F-statistics = 7.6258 and the p-value = 0.0009052 which is significant at 1%.

Hence LSAT is not significant, GPA is significant and the jointly the model is significant.

3.Test whether the size of the entering class (clsize) and the size of the faculty (faculty) need to be added to this equation jointly

Soln:

R-code:

```
modelc <- lm(log(salary)~ LSAT + GPA +log(libvol) +log(cost) +rank+ clsize+faculty,data=lawsch2)
summary(modelc)
modeld <- lm(log(salary)~ LSAT + GPA +log(libvol) +log(cost) +rank,data=lawsch2)
summary(modeld)
anova(modelc,modeld)
```

By adding the class size and faculty size we can see slight increase in R-squared value. However, the p-value for the equation jointly is 0.1381 and F-statistics is 2.0282 which show the model is insignificant.

4. What factors might influence the rank of the law school that are not included in the salary regression?

Soln:

R-code:

```
lawsch[,.(salary,rank,east,west,north,south)]
modele <- lm(log(salary)~ rank+ north+south+east+west,data=lawsch)
tidy(modele)
```

Factors like the location of school i.e., east, west, north and south aren't significant, LSAT and GPA scores are good controlling variables for rank, if we have any other factors like gender which might control the rank, however these details are not given in dataset

Question 3

$\ln[\text{price}] = \beta_0 + \beta_1 \text{sqrft} + \beta_2 \text{bdrms} + u$

$\ln[\text{price}] = 4.77 + 0.000379 \text{sqrft} + 0.0289 \text{bdrms} + u$, $n=88$, $R\text{-squared}=0.5883$

1. You are interested in estimating and obtaining a confidence interval for the percentage change in price when a 150-square-foot bedroom is added to a house. In decimal form, this is $\theta_1 = 150\beta_1 + \beta_2$. Use the data in the hprice1 table to estimate θ_1

Soln:

R-code:

```
model1 <- lm(log(price)~I(sqrft-(150*bdrms))+bdrms,data=hprice)
tidy(model1,conf.int=TRUE)
```

$\beta_0 + \beta_1 \text{sqrft} + \beta_2 \text{bdrms} + u$, $\theta_1 = 150\beta_1 + \beta_2 \rightarrow \beta_2 = \theta_1 - 150\beta_1$

$\ln[\text{price}] = \beta_0 + \beta_1 \text{sqrft} + \beta_2 \text{bdrms} + u$

$\ln[\text{price}] = \beta_0 + \beta_1 \text{sqrft} + (\theta_1 - 150\beta_1) \text{bdrms} + u$

$\ln[\text{price}] = \beta_0 + \beta_1 \text{sqrft} + (\theta_1 \text{bdrms} - 150\beta_1 \text{bdrms}) + u$

$\ln[\text{price}] = \beta_0 + \beta_1 (\text{sqrft} - 150 \text{bdrms}) + \theta_1 \text{bdrms} + u$

$\ln[\text{price}] = 4.77 + 0.000379 (\text{sqrft} - 150 \text{bdrms}) + 0.0858 \text{bdrms} + u$, $n=88$, $R\text{-squared}=0.5883$

$\theta_1 = 0.0858$ with $t\text{-statistics} = 3.21$ and $p\text{-value} = 0.00190$, so the coefficient of adding 150 sqrft bedroom is statistically significant at 1% level. Confidence interval being low 0.0326 and high 0.139

2. Write β_2 in terms of θ_1 and β_1 and plug this into the $\ln[\text{price}]$ equation.

Soln:

R-code :

```
model2 <- lm(log(price)~I(sqrft-bdrms)+bdrms,data=hprice)
tidy(model2)
summary(model2)
```

$$\theta_1 = \beta_1 + \beta_2 \rightarrow \beta_2 = \theta_1 - \beta_1$$

$$\ln[\text{price}] = \beta_0 + \beta_1 \text{sqrft} + \beta_2 \text{bdrms} + u$$

$$\ln[\text{price}] = \beta_0 + \beta_1 \text{sqrft} + (\theta_1 - \beta_1) \text{bdrms} + u$$

$$\ln[\text{price}] = \beta_0 + \beta_1 \text{sqrft} + (\theta_1 \text{bdrms} - \beta_1 \text{bdrms}) + u$$

$$\ln[\text{price}] = \beta_0 + \beta_1 (\text{sqrft} - \text{bdrms}) + \theta_1 \text{bdrms} + u$$

$$\ln[\text{price}] = 4.77 + 0.000379(\text{sqrft} - \text{bdrms}) + 0.0293 \text{bdrms} + u, n = 88 \text{ and } R\text{-squared} = 0.5883$$

3. Use part 2 to obtain a standard error for $\hat{\theta}_1$ and use this standard error to construct a 95% confidence interval.

Soln: Standard error for $\hat{\theta}_1 = 0.0296$ and 95% confidence interval is between -0.0296 and 0.0882

Question 4

$$\ln[\text{wage}] = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{tenure} + u$$

$$\ln[\text{wage}] = 5.50 + 0.0749 \text{educ} + 0.0153 \text{exper} + 0.0134 \text{tenure} + u$$

1. Consider the standard wage equation $\ln[\text{wage}] = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{tenure} + u$.

State the null hypothesis that another year of general workforce experience has the same effect on $\ln[\text{wage}]$ as another year of tenure with the current employer

Soln: Null hypothesis $\rightarrow H_0: \beta_2 = \beta_3$

2. Test the null hypothesis in part 1 against a two-sided alternative, at the 5% significance level

Soln:

R-code :

```
modela <- lm(log(wage)~educ+exper+l(exper+tenure),data=wage2)
```

```
tidy(modela,conf.int = T)
```

$$\beta_2 = \beta_3 \rightarrow \beta_2 - \beta_3 = 0 \rightarrow \beta_2 - \beta_3 = \theta_1 \rightarrow \beta_2 = \theta_1 + \beta_3$$

$$\ln[\text{wage}] = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{tenure} + u$$

$$\ln[\text{wage}] = \beta_0 + \beta_1 \text{educ} + (\theta_1 + \beta_3) \text{exper} + \beta_3 \text{tenure} + u$$

$$\ln[\text{wage}] = \beta_0 + \beta_1 \text{educ} + \theta_1 \text{exper} + \beta_3 \text{exper} + \beta_3 \text{tenure} + u$$

$$\ln[\text{wage}] = \beta_0 + \beta_1 \text{educ} + \theta_1 \text{exper} + \beta_3 (\text{exper} + \text{tenure}) + u$$

$$\ln[\text{wage}] = 5.50 + 0.0749 \text{educ} + 0.00195 \text{exper} + 0.0134 (\text{exper} + \text{tenure}) + u$$

t critical value for two-sided alternative is 0.975

The coefficient of expenditure has t-statistics = 0.412 and p-value = 0.681 is not significant at 5% level.

Hence, we cannot reject the null hypothesis in part 1

Question 5

The table 401Ksubs contains information on net financial wealth (nettfa) age of the survey respondent (age), annual family income (inc), family size (fsize), and participation in certain pension plans for people in the United States. The wealth and income variables are both recorded in thousands of dollars. For this question, use only the data for single-person households (so fsize = 1).

1. How many single-person households are there in the data set?

Soln:

R-code:

```
table(subs[,.(fsize==1)])
```

There are 2017 single person households in dataset

2. Use OLS to estimate the model $\text{nettfa} = \beta_0 + \beta_1 \text{inc} + \beta_2 \text{age} + u$, and report the results using the usual format. be sure to use only the single-person households in the sample. Interpret the slope coefficients. Are there any surprises in the slope estimates?

Soln:

R-code:

```
subs1 <- subs %>% filter(fsize == 1)
model <- lm(nettfa ~ inc + age, data = subs1)
summary(model)
augment(model)
subs1$newy <- log(augment(model)$resid^2)
model2 <- lm(newy ~ inc + age, data = subs1)
summary(model2)
subs1$h <- exp(augment(model2)$fitted)
model3 <- lm(nettfa ~ inc + age, weights = 1/h, data = subs1)
tidyw(model3)
```

$$\text{nettfa} = \beta_0 + \beta_1 \text{inc} + \beta_2 \text{age} + u$$

$$\text{nettfa} = -43 + 0.799 \text{inc} + 0.843 \text{age} + u, n=2017 \text{ and } R\text{-squared}=0.12$$

When there is an increase of one dollar in the income then the net financial wealth will increase by 80 more cents and when there is an increase in age by 1 year the net financial wealth will increase by 843\$. There is no surprise.

3. Does the intercept from the regression in part 2 have an interesting meaning? Explain

Soln:

R-code:

```
min(subs1$age)
min(subs1$inc)
```

The intercepts do not have interesting meaning here it gets the net financial wealth when the $\text{inc}=0$ and $\text{age}=0$, considering the dataset we cannot find anyone having these values.

4. Find the p-value for the test $H_0 : \beta_2 = 1$ against $H_1 : \beta_2 < 1$. Do you reject H_0 at the 1% significance level?

Soln:

R-code:

```
tstat <- (0.84266-1)/0.09202
tstat
tstat <- (0.84266-0.5)/0.09202
```

```
tstat  
pnorm(-abs(tstat))*2
```

$H_0: \beta_2 = 1 \rightarrow$ the model in part 2. t-stats=-1.709846 and p-value = 0.0873 which is not significant at 1% level hence we keep the null hypothesis and reject the alternative hypothesis H_1 .

5.If you do a simple regression of nettfa on inc, is the estimated coefficient on inc much different from the estimate in part 2? Why or why not?

Soln:

R-code:

```
model1 <- lm(nettfa~inc,data=subs1)  
summary(model1)  
cor(subs1$inc,subs1$age)
```

$\text{nettfa} = -10.571 + 0.821 \text{ inc} + u$, $n=2017$, $R\text{-squared}=0.08267$

There is a slight increase in the estimated coefficient value on inc. the difference is 0.0331. the correlation between age and income is 0.039 hence the difference.

Question 6

Use the data in the kielmc table, only for the year 1981, to answer the following questions. The data are for houses that sold during 1981 in North Andover, Massachusetts; 1981 was the year construction began on a local garbage incinerator.

1.To study the effects of the incinerator location on housing price, consider the simple regression model $\ln[\text{price}] = \beta_0 + \beta_1 \ln[\text{dist}] + u$, where price is housing price in dollars and dist is distance from the house to the incinerator measured in feet. Interpreting this equation causally, what sign do you expect for β_1 if the presence of the incinerator depresses housing prices? Estimate this equation and interpret the results.

Soln:

R-code:

```
kielmc <- kielmc %>% filter(year==1981)  
ln[price] =  $\beta_0 + \beta_1 \ln[\text{dist}] + u$   
model <- lm(log(price)~lndist,data=kielmc)  
tidy(model)  
summary(model)
```

$\ln[\text{price}] = \beta_0 + \beta_1 \ln[\text{dist}] + u$

$\ln[\text{price}] = 8.05 + 0.365 \ln[\text{dist}] + u$, $n=142$ and $R\text{-squared}=0.1803$

If there is an increase in 1% of dist then there is 36.5% increase in price of the house. Which means the effect of dist on price i.e., $\beta_1 \geq 0$

2.To the simple regression model in part 1, add the variables $\ln[\text{intst}]$, $\ln[\text{area}]$, $\ln[\text{land}]$, rooms, bath, and age, where intst is distance from the home to the interstate, area is square footage of the house, land is the lot size in square feet, rooms is total number of rooms, baths is number of bath rooms, and

age is age of the house in years. Now, what do you conclude about the effects of the incinerator? Explain why 1 and 2 give conflicting results.

Soln:

R-code:

```
model2 <- lm(log(price)~log(dist)+log(intst)+log(area)+ log(land)+ rooms+ baths+ age,data=kielmc)
summary(model2)
```

$\ln[\text{price}] = \beta_0 + \beta_1 \ln[\text{dist}] + \beta_2 \ln[\text{intst}] + \beta_3 \ln[\text{area}] + \beta_4 \ln[\text{land}] + \beta_5 \text{rooms} + \beta_6 \text{baths} + \beta_7 \text{age} + u$
 $\ln[\text{price}] = 7.59 + 0.055 \ln[\text{dist}] - 0.039 \ln[\text{intst}] + 0.32 \ln[\text{area}] + 0.076 \ln[\text{land}] + 0.043 \text{rooms} + 0.166 \text{baths} - 0.0035 \text{age} + u, n=142, R\text{-squared} = 0.7475$

we can see reduction in the value of dist from part 1 model compared to part 2 i.e., from 0.365 to 0.055 this is because of addition of other controlling variables like the number of bedrooms the area etc.

3. Add $\ln[\text{intst}]^2$ to the model from part 2. Now what happens? What do you conclude about the importance of functional form?

Soln:

R-code:

```
model3 <- lm(log(price)~log(dist)+log(intst)+log(area)+ log(land)+ rooms+ baths+ age +
l(log(intst)^2),data=kielmc)
summary(model3)
```

$\ln[\text{price}] = \beta_0 + \beta_1 \ln[\text{dist}] + \beta_2 \ln[\text{intst}] + \beta_3 \ln[\text{area}] + \beta_4 \ln[\text{land}] + \beta_5 \text{rooms} + \beta_6 \text{baths} + \beta_7 \text{age} + \beta_8 \ln[\text{intst}]^2 + u$

$\ln[\text{price}] = -3.318025 + 0.185256 \ln[\text{dist}] + 2.072959 \ln[\text{intst}] + 0.359352 \ln[\text{area}] + 0.091386 \ln[\text{land}] + 0.038106 \text{rooms} + 0.149533 \text{baths} - 0.002927 \text{age} + -0.119329 \ln[\text{intst}]^2 + u$

After adding the $\ln[\text{intst}]^2$ the distance from house to incinerator has become slightly significant, the distance to interstate inst has also become significant along with rooms and bathrooms. We can also see this increase in significance with dist and intst as there is a strong positive correlation between these two factors.

4. Is the square of $\ln[\text{dist}]$ significant when you add it to the model from part 3?

Soln:

R-code:

```
model4 <- lm(log(price)~log(dist)+log(intst)+log(area)+ log(land)+ rooms+ baths+ age +
l(log(intst)^2)+l(ldist^2),data=kielmc)
summary(model4)
```

Adding $\ln[\text{dist}]^2$ doesn't make the model better as the coefficient is not significant even at 10% level.

Question 7

Use the data in the wage1 table for this exercise.

1. Use OLS to estimate the equation $\ln[\text{wage}] = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{exper}^2 + u$ and report the results using the usual format

Soln:

R-code:

```
model <- lm(log(wage)~educ+exper+l(exper^2),data=wage1)
summary(model)
```

$$\ln[\text{wage}] = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{exper}^2 + u$$

$$\ln[\text{wage}] = 0.126 + 0.0906 \text{educ} + 0.0410 \text{exper} - 0.000712 \text{exper}^2 + u, n=526, R\text{-squared}=0.3003$$

2. Is exper² statistically significant at the 1% level?

Soln: The coefficients of exper^2 with t-stats = -6.141 and p-value = 0.0000000016276762 is significant at 1% level

3. Using the approximation $\partial \ln[\text{wage}] / \partial \text{exper} \approx \beta^2 + 2\beta^3 \text{exper}$, find the approximate return to the fifth year of experience. What is the approximate return to the twentieth year of experience?

Soln: for the return to 5th year of experience $\partial \ln[\text{wage}] \approx 0.0410 + 2 * (-0.000712) * 4 = 0.0035304$

For the return to 20th year of experience $\partial \ln[\text{wage}] \approx 0.0410 + 2 * (-0.000712) * 19 = 0.013944$

4. At what value of exper does additional experience actually lower predicted $\ln[\text{wage}]$? How many people have more experience in this sample?

Soln:

R-code:

```
count_exp <- wage1 %>% filter(exper >= 29)
nrow(count_exp)
```

lower predicted $\ln[\text{wage}] \approx \beta^2 + 2\beta^3 \text{exper} \approx 28.79 \approx 29$ years i.e., people in data set with at least 29 years of experience will lower the predicted $\ln[\text{wage}]$ we have 121 people with greater than or equal to 29 years of experience.

Question 8

Consider a model where the return to education depends upon the amount of work experience (and vice versa): $\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{educ} \times \text{exper} + u$

$$\ln[\text{wage}] = 5.95 + 0.0440 \text{educ} - 0.0215 \text{exper} + 0.00320 (\text{educ} * \text{exper}) + u$$

1. Show that the return to another year of education (in decimal form), holding experience fixed, is $\beta_1 + \beta_3 \text{exper}$

Soln: By holding experience fixed above equation will be

$$\partial \ln[\text{wage}] = \beta_1 \partial(\text{educ}) + \beta_3 \partial(\text{educ}) * \text{exper}$$

$$\partial \ln[\text{wage}] = \partial(\text{educ}) (\beta_1 + \beta_3 * \text{exper})$$

$$\partial \ln[\text{wage}] / \partial(\text{educ}) = \beta_1 + \beta_3 * \text{exper}$$

2. State the null hypothesis that the return to education does not depend on the level of experience. What do you think is the appropriate alternative?

Soln:

R-code:

$H_0: \beta_3=0$ is the null hypothesis and $H_1: \beta_3>0$ is the alternative hypothesis which states return to education depends on the level of experience

3. Use the data in the wage2 table to test the null hypothesis in 2 against your stated alternative

Soln:

R-code:

```
model <- lm(log(wage)~educ+exper,data=wage2)
tidy(model)
```

5.95+0.0440 educ-0.0215exper+0.00320(exper*educ) +u

t-stats =2.09 and p-value= 0.0365 of the interaction term is not significant even at 1% level hence we reject the H_0 against the H_1

4. Let θ_1 denote the return to education (in decimal form), when $\text{exper} = 10$: $\theta_1 = \beta_1 + 10\beta_3$. Estimate θ_1 and a 95% confidence interval for θ_1 . (Hint: Write $\beta_1 = \theta_1 - 10\beta_3$ and plug this into the equation; then rearrange. This gives the regression for obtaining the confidence interval for θ_1 .)

Soln:

R-code:

```
modelb <- lm(log(wage)~educ+exper+l(educ*(-10+exper)),data=wage2)
tidy(modelb,conf.int = TRUE)
```

$$\beta_1 = \theta_1 - 10\beta_3$$

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{educ} \times \text{exper} + u$$

$$\log(\text{wage}) = \beta_0 + (\theta_1 - 10\beta_3) \text{educ} + \beta_2 \text{exper} + \beta_3 \text{educ} \times \text{exper} + u$$

$$\log(\text{wage}) = \beta_0 + \theta_1 \text{educ} - 10\beta_3 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{educ} \times \text{exper} + u$$

$$\log(\text{wage}) = \beta_0 + \theta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{educ} (-10 + \text{exper}) + u$$

$$5.95+0.761 \text{educ}-0.0215 \text{exper}+0.00320 \text{educ}(\text{exper}-10) + u$$

$\theta_1 = 0.761$ with 95% confidence interval between 0.0631 and 0.0891

Question 9

Use the data in the gpa2 table for this exercise.

1. Estimate the model $\text{sat} = \beta_0 + \beta_1 \text{hsize} + \beta_2 \text{hsize}^2 + u$. where hsize is the size of the graduating class (in hundreds) and write the results in the usual form. Is the quadratic term statistically significant?

Soln:

R-code:

```
model <- lm(sat~hsize+l(hsize^2),data=gpa2)
tidy(model)
summary(model)
```

$$\text{sat} = \beta_0 + \beta_1 \text{hsize} + \beta_2 \text{hsize}^2 + u$$

$$\text{sat} = 998 + 19.8 \text{hsize} - 2.13 \text{hsize}^2 + u, n=935, R\text{-squared}=0.1551$$

the coefficient of quadratic term is statistically significant at 1% level with t-stat = -3.90 and p-value=0.0000960

2. Using the estimated equation from part 1, what is the “optimal” high school size? Justify your answer

Soln: $y = -b/2a$ which is the turning point in a parabola, $-19.8/(2 \cdot -2.13) = 19.8/4.26 = 4.65$ i.e., 465 will be the optimal high school size

3. Is this analysis representative of the academic performance of all high school seniors? Explain.

Soln: No this doesn't represent academic performance of all high school seniors as the data set does not specify the SAT scores of high school student it's just the collective combined score of student who took SAT exam.

4. Find the estimated optimal high school size, using $\ln(\text{sat})$ as the dependent variable. Is it much different from what you obtained in part 2?

Soln:

R-code:

```
modela <- lm(log(sat)~hsize+l(hsize^2),data=gpa2)
tidy(modela)
```

$\ln(\text{sat}) = 6.9 + 0.0196 \text{hsize} - 0.00209 \text{Hsize}^2 + u$ here the optimal high school size is $0.0196/(2 \cdot 0.00209) = 0.0196/0.00418 = 4.68$ i.e., 468 which is close to the value obtained in part 2

Question 10

Use the housing price data in the hprice1 table for this exercise.

1. Estimate the model $\ln[\text{price}] = \beta_0 + \beta_1 \ln[\text{lotsize}] + \beta_2 \ln[\text{sqrft}] + \beta_3 \text{bdrms} + u$ and report the results in the usual OLS format

Soln:

R-code:

```
summary(lm(log(price)~log(lotsize)+log(sqrft)+bdrms,data=hprice))
```

$$\ln[\text{price}] = \beta_0 + \beta_1 \ln[\text{lotsize}] + \beta_2 \ln[\text{sqrft}] + \beta_3 \text{bdrms} + u$$

$$\ln[\text{price}] = -1.29 + 0.168 \ln(\text{lotsize}) + 0.7 \ln(\text{sqrft}) + 0.037 \ln(\text{bdrms}) + u, n=88 R\text{-squared} = 0.643$$

2. Find the predicted value of price, when lotsize = 20 000, sqrft = 2 500, and bdrms = 4

Soln:

R-code:

```
model <- lm(log(price)~log(lotsize)+log(sqrft)+bdrms,data=hprice)
predict_model <- predict(model,data.frame(lotsize = 20000,sqrft = 2500,bdrms = 4))
exp(predict_model)
```

Predicted value = 400.574 → 400574\$

3. For explaining variation in price, decide whether you prefer the model from part 1 or the model
 $\text{price} = \beta_0 + \beta_1 \text{lotsize} + \beta_2 \text{sqrft} + \beta_3 \text{bdrms} + u$

Soln:

R-code:

```
summary(lm(price~lotsize+sqrft+bdrms,data=hprice,scipen=999))
```

$\text{price} = -21.77 + 0.00206 \text{lotsize} + 0.1227 \text{sqrft} + 13.85 \text{bdrms} + u$, $n=88$, $R\text{-squared}=0.6724$

For explaining variation in price we prefer model in part 1 which is easy to interpret as the dependent variables continuous type and log is better.