

Lecture 18 - ACS Example

Sta102 / BME 102

Colin Rundel

April 9, 2013

1 Transformations and ACS

- Diagnostics
- Variance stabilizing transformations
- New model: log of income
- Diagnostics for model for logged income model
- Interpretations for model for logged income

From last lab

Just like in lab we load data, and subset for those who were employed.

```
download.file("http://stat.duke.edu/~cr173/Sta102_Sp14/Lab/acs.RData",  
             destfile = "acs.RData")  
load("acs.RData")  
acs_sub = subset(acs, acs$employment == "employed")
```

Predicting income

```
l = lm(income ~ hrs_work + race + age + gender + edu + disability, data = acs_sub)
summary(l)
```

```
##
## Call:
## lm(formula = income ~ hrs_work + race + age + gender + edu +
##     disability, data = acs_sub)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -122650  -20503   -4597   10945   321681
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -21737       7719   -2.82  0.00498 **
## hrs_work         1000         136    7.38  3.9e-13 ***
## raceblack      -6016        5877   -1.02  0.30636
## raceasian     29596        8030    3.69  0.00024 ***
## raceother     -8599        6649   -1.29  0.19624
## age             562         119    4.72  2.7e-06 ***
## genderfemale  -18121       3496   -5.18  2.7e-07 ***
## educollege     17274       3828    4.51  7.3e-06 ***
## edugrad        58552       5419   10.81 < 2e-16 ***
## disabilityyes  -15852       6210   -2.55  0.01086 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 48300 on 833 degrees of freedom
## Multiple R-squared:  0.294, Adjusted R-squared:  0.286
## F-statistic: 38.5 on 9 and 833 DF,  p-value: <2e-16
```

Categorical variables with multiple levels

In model selection based on p-value:

Leave variable in the model if p-value for *any* level is below α_{crit} .

For example, the race variable in our model:

	Estimate	Std. Error	t value	Pr(> t)
...				
raceblack	-6015.53	5877.30	-1.02	0.31
raceasian	29595.59	8029.98	3.69	0.00
raceother	-8599.21	6648.63	-1.29	0.20
...				

Categorical variables with multiple levels

In model selection based on p-value:

Leave variable in the model if p-value for *any* level is below α_{crit} .

For example, the race variable in our model:

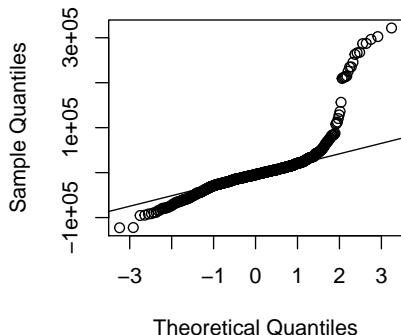
	Estimate	Std. Error	t value	Pr(> t)
...				
raceblack	-6015.53	5877.30	-1.02	0.31
raceasian	29595.59	8029.98	3.69	0.00
raceother	-8599.21	6648.63	-1.29	0.20
...				

How do we interpret the slopes associated with the race variable?

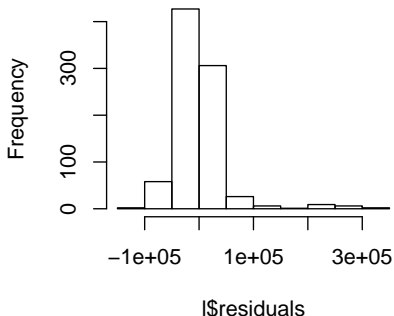
(1) Nearly normal residuals

```
par(mfrow = c(1, 2))  
qqnorm(l$residuals, main = "Normal prob. plot\nof residuals")  
qqline(l$residuals)  
hist(l$residuals, main = "Histogram of residuals")
```

**Normal prob. plot
of residuals**



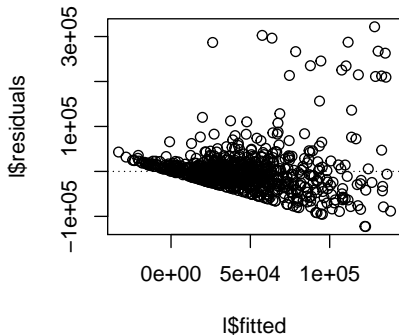
Histogram of residuals



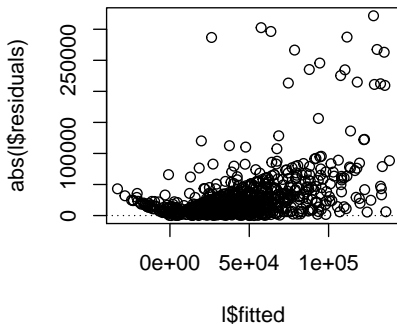
(2) Constant variability of residuals

```
par(mfrow = c(1, 2))  
plot(l$residuals ~ l$fitted, main = "Residuals vs. fitted")  
abline(h = 0, lty = 3)  
plot(abs(l$residuals) ~ l$fitted, main = "Absolute value of\nresiduals vs. fitted")  
abline(h = 0, lty = 3)
```

Residuals vs. fitted



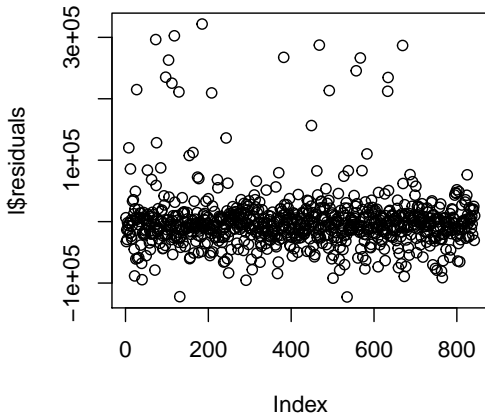
Absolute value of residuals vs. fitted



(3) Independence

```
plot(l$residuals, main = "Residuals vs. order of data collection")
```

Residuals vs. order of data collection

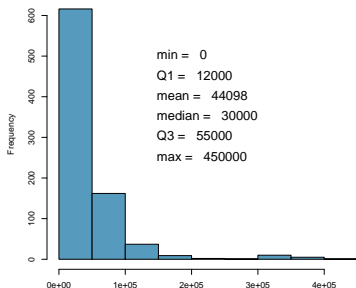


Transformations

- We saw that residuals have a right-skewed distribution, and the relationship between hours worked per week and income is non-linear (exponential).
- In these situations a transformation applied to the response variable may be useful.
- In order to decide which transformation to use, we should examine the distribution of the response variable.

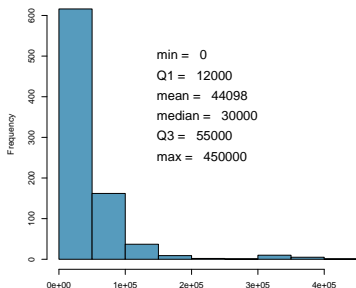
Transformations

- We saw that residuals have a right-skewed distribution, and the relationship between hours worked per week and income is non-linear (exponential).
- In these situations a transformation applied to the response variable may be useful.
- In order to decide which transformation to use, we should examine the distribution of the response variable.



Transformations

- We saw that residuals have a right-skewed distribution, and the relationship between hours worked per week and income is non-linear (exponential).
- In these situations a transformation applied to the response variable may be useful.
- In order to decide which transformation to use, we should examine the distribution of the response variable.



- The distribution is right skewed → suggests that a log transformation may be useful.

Log of 0

```
summary(acs_sub$income)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##         0  12000   30000   44100   55000  450000
```

```
log(0)
```

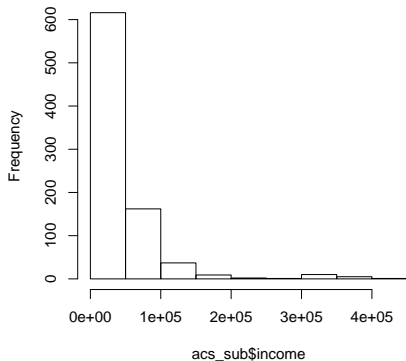
```
## [1] -Inf
```

- Since there are some individuals who had 0 income (from salaries and wages) last year, we cannot take the log of their income, since $\log(0) = -\infty$.
- A commonly used trick is to add a very small number to all values before taking the log.

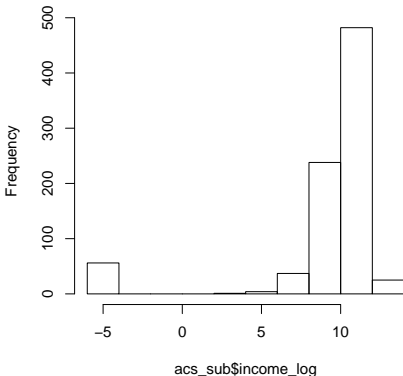
Logged income distribution

```
acs_sub$income_log = log(acs_sub$income + 0.01)
par(mfrow = c(1, 2))
hist(acs_sub$income)
hist(acs_sub$income_log)
```

Histogram of acs_sub\$income

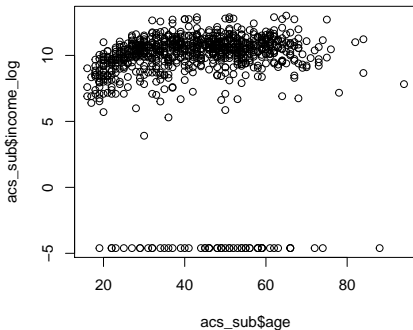
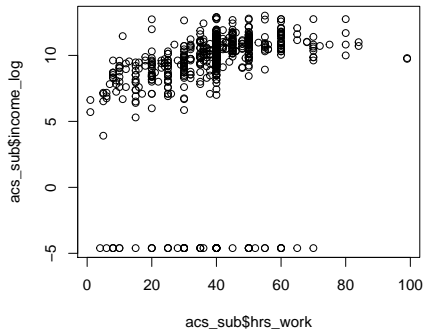


Histogram of acs_sub\$income_log



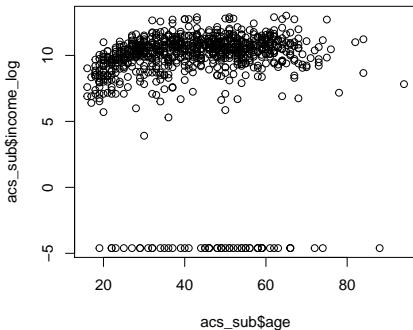
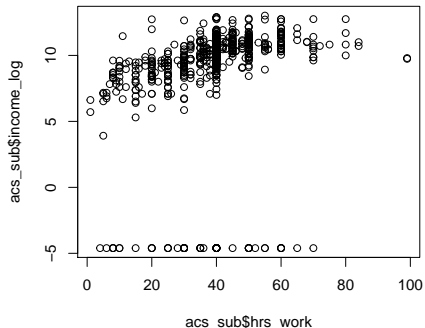
Logged income relationships

```
par(mfrow = c(1, 2), mar = c(5, 4, 1, 2) + 0.1)
plot(acs_sub$income_log ~ acs_sub$hrs_work)
plot(acs_sub$income_log ~ acs_sub$age)
```



Logged income relationships

```
par(mfrow = c(1, 2), mar = c(5, 4, 1, 2) + 0.1)
plot(acs_sub$income_log ~ acs_sub$hrs_work)
plot(acs_sub$income_log ~ acs_sub$age)
```



We still might want to do something about those 0 incomes, it doesn't make sense to model them with the rest of the data.

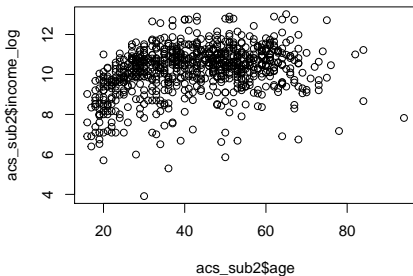
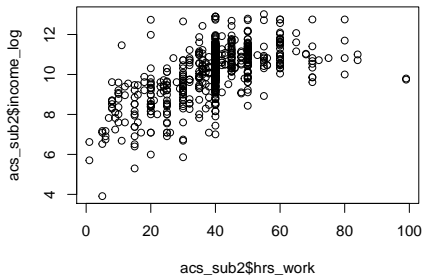
Further subsetting the data

People who work more than 0 hours per week but make 0 income in salaries and wages are different than others whose income is proportional to number of hours they work. So we have reason to omit these people from the analysis (and model their income differently based on other variables).

```
acs_sub2 = subset(acs_sub, acs_sub$income > 0)
acs_sub2$income_log = log(acs_sub2$income)
```

Logged relationships - for those with any income

```
par(mfrow = c(1, 2))  
plot(acs_sub2$income_log ~ acs_sub2$hrs_work)  
plot(acs_sub2$income_log ~ acs_sub2$age)
```



Predicting log of income

```

l_log = lm(income_log ~ hrs_work + race + age + gender + edu +
  disability, data = acs_sub2)
summary(l_log)

##
## Call:
## lm(formula = income_log ~ hrs_work + race + age + gender + edu +
##     disability, data = acs_sub2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.573 -0.394  0.088  0.499  3.165
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   7.31368    0.14074   51.97 < 2e-16 ***
## hrs_work      0.04879    0.00253   19.32 < 2e-16 ***
## raceblack    -0.14758    0.10436   -1.41  0.16
## raceasian    0.13688    0.14162    0.97  0.33
## raceother   -0.19219    0.12119   -1.59  0.11
## age          0.02223    0.00217   10.22 < 2e-16 ***
## genderfemale -0.27608    0.06370   -4.33 1.7e-05 ***
## educollege   0.39923    0.06993    5.71 1.6e-08 ***
## edugrad      0.83369    0.09871    8.45 < 2e-16 ***
## disabilityyes -0.62448    0.11549   -5.41 8.5e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.849 on 777 degrees of freedom
## Multiple R-squared:  0.52, Adjusted R-squared:  0.515
## F-statistic: 93.6 on 9 and 777 DF,  p-value: <2e-16

```

Final model for log of income

```

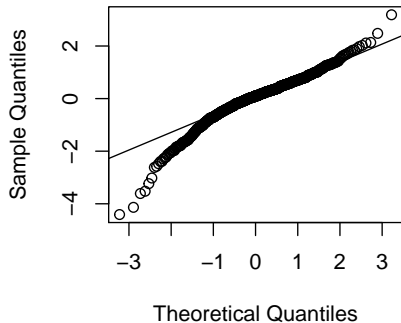
l_log_final = lm(income_log ~ hrs_work + age + gender + edu +
  disability, data = acs_sub2)
summary(l_log_final)

##
## Call:
## lm(formula = income_log ~ hrs_work + age + gender + edu + disability,
##     data = acs_sub2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.410 -0.394  0.099  0.512  3.188
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   7.28126    0.13908   52.35 < 2e-16 ***
## hrs_work       0.04902    0.00253   19.39 < 2e-16 ***
## age            0.02231    0.00217   10.30 < 2e-16 ***
## genderfemale  -0.28736    0.06357   -4.52 7.1e-06 ***
## educollege     0.41356    0.06974    5.93 4.6e-09 ***
## edugrad        0.84491    0.09832    8.59 < 2e-16 ***
## disabilityyes -0.63204    0.11556   -5.47 6.1e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.85 on 780 degrees of freedom
## Multiple R-squared:  0.517, Adjusted R-squared:  0.513
## F-statistic: 139 on 6 and 780 DF,  p-value: <2e-16

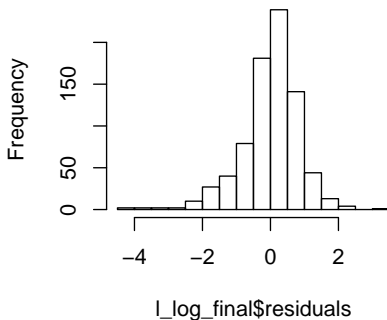
```

(1) Nearly normal residuals

Normal prob. plot of residuals

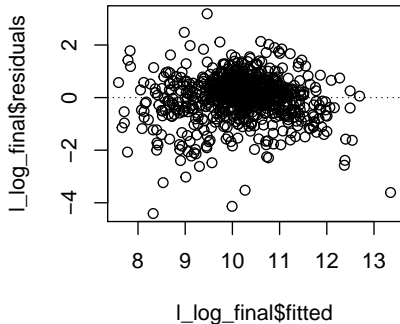


Histogram of residuals

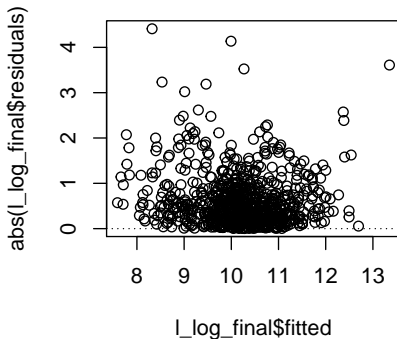


(2) Constant variability of residuals

Residuals vs. fitted



Absolute value of residuals vs. fitted



Residuals vs. fitted

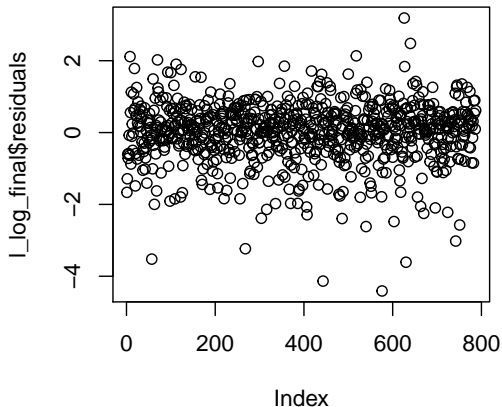


Absolute value of residuals vs. fitted



(3) Independence

Residuals vs. order of data collection



Interpretation

Which of the following is the correct interpretation of the slope of age hours worked per week?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	7.28	0.14	52.35	0.00
hrs_work	0.05	0.00	19.39	0.00
age	0.02	0.00	10.30	0.00
genderfemale	-0.29	0.06	-4.52	0.00
educollege	0.41	0.07	5.93	0.00
edugrad	0.84	0.10	8.59	0.00
disabilityyes	-0.63	0.12	-5.47	0.00

Interpretation

Which of the following is the correct interpretation of the slope of age hours worked per week?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	7.28	0.14	52.35	0.00
hrs_work	0.05	0.00	19.39	0.00
age	0.02	0.00	10.30	0.00
genderfemale	-0.29	0.06	-4.52	0.00
educollege	0.41	0.07	5.93	0.00
edugrad	0.84	0.10	8.59	0.00
disabilityyes	-0.63	0.12	-5.47	0.00

For each additional hour worked per week, we would expect income to increase on average by a factor of 105.12% ($\exp(0.05) = 1.0512$).

Interpretation (cont.)

Which of the following is the correct interpretation of the slope of edu:college?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	7.28	0.14	52.35	0.00
hrs_work	0.05	0.00	19.39	0.00
age	0.02	0.00	10.30	0.00
gender:female	-0.29	0.06	-4.52	0.00
edu:college	0.41	0.07	5.93	0.00
edu:grad	0.84	0.10	8.59	0.00
disability:yes	-0.63	0.12	-5.47	0.00

Interpretation (cont.)

Which of the following is the correct interpretation of the slope of edu:college?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	7.28	0.14	52.35	0.00
hrs_work	0.05	0.00	19.39	0.00
age	0.02	0.00	10.30	0.00
gender:female	-0.29	0.06	-4.52	0.00
edu:college	0.41	0.07	5.93	0.00
edu:grad	0.84	0.10	8.59	0.00
disability:yes	-0.63	0.12	-5.47	0.00

The model predicts that college educated individuals, on average, make 150.7% more ($\exp(0.41) = 1.507$) than those who have a HS degree or lower.