

Lab #7 - More on Regression in R

Econ 224

September 18th, 2018

Robust Standard Errors

Your reading assignment from Chapter 3 of ISL briefly discussed two ways that the standard regression inference formulas built into R can go wrong: (1) non-constant error variance, and (2) correlation between regression errors. Today we'll briefly look at the first of these problems and how to correct for it.

Consider the simple linear regression $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$. If the variance of ϵ_i is unrelated to the value of the predictor x_i , we say that the regression errors are *homoskedastic*. This is just a fancy Greek word for *constant variance*. If instead, the variance of ϵ_i depends on the value of x_i , we say that the regression errors are *heteroskedastic*. This is just a fancy Greek word for *non-constant variance*. Heteroskedasticity does not invalidate our least squares estimates of β_0 and β_1 , but it does invalidate the formulas used by `lm` to calculate standard errors and p-values.

Let's look at a simple simulation example:

```
set.seed(4321)
n <- 100
x <- runif(n)
e1 <- rnorm(n, mean = 0, sd = sqrt(2 * x))
e2 <- rnorm(n, mean = 0, sd = 1)
intercept <- 0.2
slope <- 0.9
y1 <- intercept + slope * x + e1
y2 <- intercept + slope * x + e2
library(tidyverse)
mydat <- tibble(x, y1, y2)
rm(x, y1, y2)
```

Exercise #1

[approx 10 min]

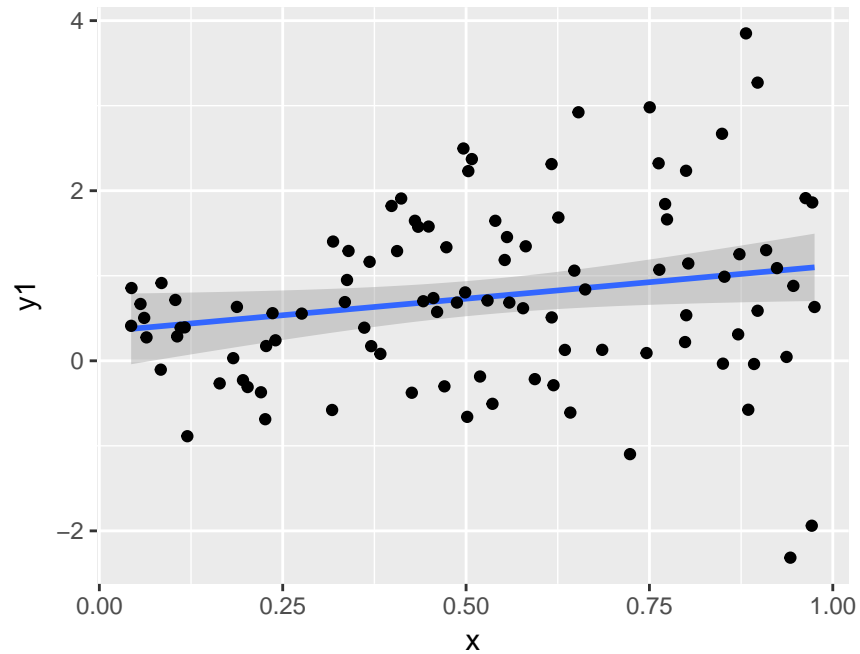
1. Read through my simulation code and make sure you understand what each step is going. What is the distribution of the errors? What is the distribution of x ? In the simulation design, is there a relationship between x and $y1$? What about $y2$?
2. For each of the two simulated outcome variables $y1$ and $y2$, plot the outcome against x along with the linear regression line.
3. Based on your plots from part 2 and the simulation code, which errors are heteroskedastic: $e1$, $e2$, both, or neither? How can you tell?

Solution to Exercise #1

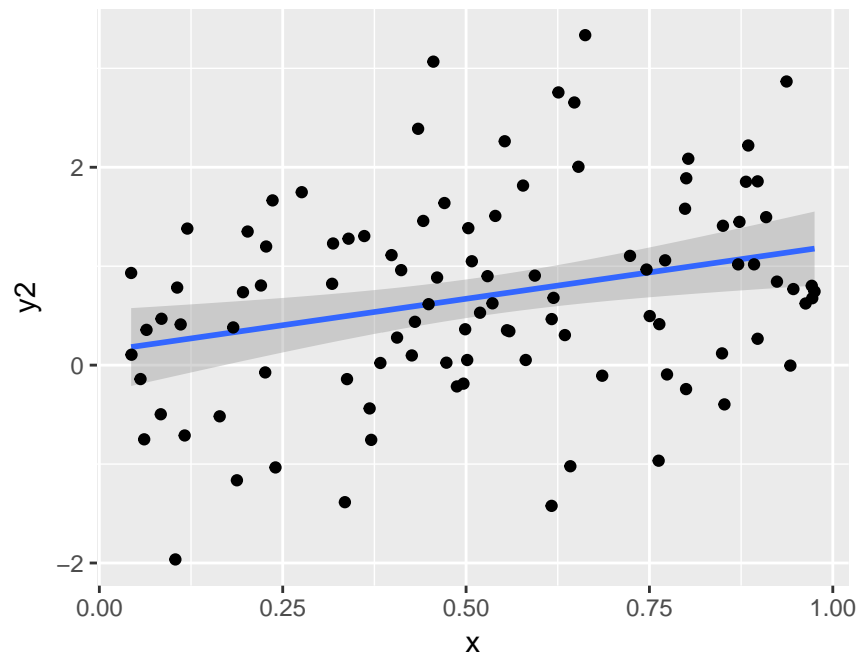
1. x is uniform and the errors are normally distributed. There is indeed a relationship between x and y : the conditional mean of $y1$ given x is $0.2 + 0.4 x$ and the same is true of $y2$

2. Here is a simple way to make the plots:

```
library(ggplot2)
ggplot(mydat, aes(x, y1)) +
  geom_smooth(method = 'lm') +
  geom_point()
```



```
ggplot(mydat, aes(x, y2)) +
  geom_smooth(method = 'lm') +
  geom_point()
```



3. The errors `e1` are heteroskedastic while the errors `e2` are homoskedastic. We can see this both from plotting the data which “fan out” around the regression line for `y1` and from the simulation code: to generate `e1` we multiplied some normal random draws by the value of `x` so the variance clearly depends on `x`

Robust Standard Errors using `lm_robust`

Install the package `estimatr`. Provides a replacement for `lm` called `lm_robust` that allows us to choose robust standard errors

```
library(estimatr)
reg1_classical <- lm_robust(y1 ~ x, mydat, se_type = 'stata')
summary(reg1_classical)
```

Call:

```
lm_robust(formula = y1 ~ x, data = mydat, se_type = "stata")
```

Standard error type: HC1

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	CI Lower	CI Upper	DF
(Intercept)	0.3418	0.1739	1.966	0.05215	-0.003241	0.6868	98
x	0.7766	0.4068	1.909	0.05919	-0.030707	1.5839	98

Multiple R-squared: 0.04119 , Adjusted R-squared: 0.0314

F-statistic: 3.644 on 1 and 98 DF, p-value: 0.05919

```
reg1_robust <- lm_robust(y1 ~ x, mydat, se_type = 'classical')
summary(reg1_robust)
```

Call:

```
lm_robust(formula = y1 ~ x, data = mydat, se_type = "classical")
```

Standard error type: classical

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	CI Lower	CI Upper	DF
(Intercept)	0.3418	0.2240	1.526	0.13027	-0.10273	0.7863	98
x	0.7766	0.3785	2.052	0.04286	0.02548	1.5277	98

Multiple R-squared: 0.04119 , Adjusted R-squared: 0.0314

F-statistic: 4.21 on 1 and 98 DF, p-value: 0.04286

The nice thing about using `lm_robust` is that it plays nicely with `linearHypothesis` for carrying out F-tests. In an example with only one regressor the F-test is completely superfluous (the F-test statistic is simply the square of the t-test statistic for the slope!) but just to see that it works:

```
library(car)
summary(lm(y1 ~ x, mydat))$fstatistic
```

```

      value      numdf      dendif
4.209829  1.000000  98.000000

```

```
linearHypothesis(reg1_classical, 'x = 0')
```

Linear hypothesis test

Hypothesis:

$x = 0$

Model 1: restricted model

Model 2: $y1 \sim x$

```

      Res.Df Df   Chisq Pr(>Chisq)
1         99
2         98  1 3.6442    0.05626 .
---

```

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```
linearHypothesis(reg1_robust, 'x = 0')
```

Linear hypothesis test

Hypothesis:

$x = 0$

Model 1: restricted model

Model 2: $y1 \sim x$

```

      Res.Df Df   Chisq Pr(>Chisq)
1         99
2         98  1 4.2098    0.04019 *
---

```

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Exercise #2

[approx 5 min]

Repeat my inference comparison from above for the regression $y2 \sim x$ using classical and robust standard errors. Explain your results. Do we need to use robust standard errors in this case? Why or why not?

Solution to Exercise #2

We do not in fact need robust standard errors in this case since the errors are homoskedastic. If we wanted to be on the safe side, we could still use them, but notice that the results are slightly different.

```

reg2_classical <- lm_robust(y2 ~ x, mydat, se_type = 'stata')
summary(reg1_classical)

```

```
Call:
lm_robust(formula = y1 ~ x, data = mydat, se_type = "stata")
```

Standard error type: HC1

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	CI Lower	CI Upper	DF
(Intercept)	0.3418	0.1739	1.966	0.05215	-0.003241	0.6868	98
x	0.7766	0.4068	1.909	0.05919	-0.030707	1.5839	98

Multiple R-squared: 0.04119 , Adjusted R-squared: 0.0314

F-statistic: 3.644 on 1 and 98 DF, p-value: 0.05919

```
reg2_robust <- lm_robust(y2 ~ x, mydat, se_type = 'classical')
summary(reg1_robust)
```

Call:

```
lm_robust(formula = y1 ~ x, data = mydat, se_type = "classical")
```

Standard error type: classical

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	CI Lower	CI Upper	DF
(Intercept)	0.3418	0.2240	1.526	0.13027	-0.10273	0.7863	98
x	0.7766	0.3785	2.052	0.04286	0.02548	1.5277	98

Multiple R-squared: 0.04119 , Adjusted R-squared: 0.0314

F-statistic: 4.21 on 1 and 98 DF, p-value: 0.04286

```
summary(lm(y2 ~ x, mydat))$fstatistic
```

value	numdf	dendf
8.899712	1.000000	98.000000

```
linearHypothesis(reg2_classical, 'x = 0')
```

Linear hypothesis test

Hypothesis:

x = 0

Model 1: restricted model

Model 2: y2 ~ x

	Res.Df	Df	Chisq	Pr(>Chisq)
1	99			
2	98	1	10.89	0.0009669 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```
linearHypothesis(reg2_robust, 'x = 0')
```

Linear hypothesis test

Hypothesis:

$x = 0$

Model 1: restricted model

Model 2: $y2 \sim x$

	Res.Df	Df	Chisq	Pr(>Chisq)
1	99			
2	98	1	8.8997	0.002852 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Exercise

Go back and do robust standard errors for one of the questions from the previous problem set.

Publication-quality Tables with **stargazer**

We'll use the package **stargazer** to generate pretty tables of results like the ones you see in journal articles. Make sure to install this package before proceeding.

```
library(stargazer)
```

Simple table of summary statistics

I chose to output my .Rmd file to a pdf using LaTeX, so I used the option `type = latex`. If you're using `html` you'll need to change this to `type = 'html'`. If you want to see a “preview” of the table within R studio without compiling, choose `type = 'text'`. Also notice I'm using the `knitr` option `asis`. You'll need this to make sure that the **stargazer** table knits correctly.

```
stargazer(mtcars, type = 'latex')
```

```
% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu
% Date and time: Fri, Sep 07, 2018 - 05:55:43 PM
```

The **stargazer** command provides dozens of options for customizing the appearance of the output it generates. Here's a nicer version of the preceding table that uses some of these options:

```
mylabels <- c('Miles/gallon',
              'No. of cylinders',
              'Displacement (cubic inches)',
              'Horsepower',
              'Rear axle ratio',
              'Weight (1000lb)',
              '1/4 Mile Time',
```

Table 1:

Statistic	N	Mean	St. Dev.	Min	Pctl(25)	Pctl(75)	Max
mpg	32	20.091	6.027	10	15.4	22.8	34
cyl	32	6.188	1.786	4	4	8	8
disp	32	230.722	123.939	71	120.8	326	472
hp	32	146.688	68.563	52	96.5	180	335
drat	32	3.597	0.535	2.760	3.080	3.920	4.930
wt	32	3.217	0.978	1.513	2.581	3.610	5.424
qsec	32	17.849	1.787	14.500	16.892	18.900	22.900
vs	32	0.438	0.504	0	0	1	1
am	32	0.406	0.499	0	0	1	1
gear	32	3.688	0.738	3	3	4	5
carb	32	2.812	1.615	1	2	4	8

```

      'V/S',
      'Manual Transmission? (1 = Yes)',
      'No. forward gears',
      'No. carburetors')
stargazer(mtcars,
  type = 'latex',
  title = 'Summary Statistics: Motor Trend Cars Dataset',
  digits = 1,
  header = FALSE,
  covariate.labels = mylabels)

```

Table 2: Summary Statistics: Motor Trend Cars Dataset

Statistic	N	Mean	St. Dev.	Min	Pctl(25)	Pctl(75)	Max
Miles/gallon	32	20.1	6.0	10	15.4	22.8	34
No. of cylinders	32	6.2	1.8	4	4	8	8
Displacement (cubic inches)	32	230.7	123.9	71	120.8	326	472
Horsepower	32	146.7	68.6	52	96.5	180	335
Rear axle ratio	32	3.6	0.5	2.8	3.1	3.9	4.9
Weight (1000lb)	32	3.2	1.0	1.5	2.6	3.6	5.4
1/4 Mile Time	32	17.8	1.8	14.5	16.9	18.9	22.9
V/S	32	0.4	0.5	0	0	1	1
Manual Transmission? (1 = Yes)	32	0.4	0.5	0	0	1	1
No. forward gears	32	3.7	0.7	3	3	4	5
No. carburetors	32	2.8	1.6	1	2	4	8

We can also customize which summary statistics are reported using the options `summary.stat` and `omit.summary.stat`. For example, if we only wanted to show the mean, standard deviation, and quartiles of the data, we could use the following:

```

stargazer(mtcars,
  type = 'latex',
  title = 'Summary Statistics: Motor Trend Cars Dataset',
  digits = 1,
  header = FALSE,

```

```

covariate.labels = mylabels,
summary.stat = c('mean',
                 'sd',
                 'p25',
                 'median',
                 'p75'))

```

Table 3: Summary Statistics: Motor Trend Cars Dataset

Statistic	Mean	St. Dev.	Pctl(25)	Median	Pctl(75)
Miles/gallon	20.1	6.0	15.4	19.2	22.8
No. of cylinders	6.2	1.8	4	6	8
Displacement (cubic inches)	230.7	123.9	120.8	196.3	326
Horsepower	146.7	68.6	96.5	123	180
Rear axle ratio	3.6	0.5	3.1	3.7	3.9
Weight (1000lb)	3.2	1.0	2.6	3.3	3.6
1/4 Mile Time	17.8	1.8	16.9	17.7	18.9
V/S	0.4	0.5	0	0	1
Manual Transmission? (1 = Yes)	0.4	0.5	0	0	1
No. forward gears	3.7	0.7	3	4	4
No. carburetors	2.8	1.6	2	2	4

Exercise

Figure out what each of the options `title`, `digits`, `header` and `covariate.labels` does in the preceding code chunk. Write a one sentence explanation of each. Try experimenting with different values and or consulting the documentation for `stargazer`.

Exercise

Make a pretty table of summary statistics from one of the recent problem set questions.

Regression Output

If you pass a dataframe (or tibble) to `stargazer`, by default it will create a table of summary statistics. If you instead pass a *regression* object, it will make a regression table. For example: Run a bunch of regressions using `mtcars`

```

reg1 <- lm(mpg ~ disp, mtcars)
stargazer(reg1, type = 'latex',
          header = FALSE,
          digits = 1,
          title = 'Predicting Fuel Economy from Displacement')

```

Let's run a few more regressions and make a table that summarizes the results of *all* of them:

Table 4: Predicting Fuel Economy from Displacement

<i>Dependent variable:</i>	
	mpg
disp	-0.04*** (0.005)
Constant	29.6*** (1.2)
Observations	32
R ²	0.7
Adjusted R ²	0.7
Residual Std. Error	3.3 (df = 30)
F Statistic	76.5*** (df = 1; 30)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

```
reg2 <- lm(mpg ~ wt, mtcars)
reg3 <- lm(mpg ~ disp + wt, mtcars)
stargazer(reg1, reg2, reg3,
  type = 'latex',
  digits = 1,
  header = FALSE,
  title = 'Regression Results for Motor Trend Dataset',
  covariate.labels = c('Displacement (cubic inches)', 'Weight (1000lb)'),
  dep.var.caption = 'Miles/gallon',
  notes = c('Data are courtesy of Motor Trend Magazine. Also, R rules!'))
```

Notice how I added a label for the dependent variable and appended a *note* to the regression table.

Exercise

Make a pretty regression table for one of the problems on the last problem set. Use robust standard errors.

Table 5: Regression Results for Motor Trend Dataset

	Miles/gallon		
	mpg		
	(1)	(2)	(3)
Displacement (cubic inches)	-0.04*** (0.005)		-0.02* (0.01)
Weight (1000lb)		-5.3*** (0.6)	-3.4*** (1.2)
Constant	29.6*** (1.2)	37.3*** (1.9)	35.0*** (2.2)
Observations	32	32	32
R ²	0.7	0.8	0.8
Adjusted R ²	0.7	0.7	0.8
Residual Std. Error	3.3 (df = 30)	3.0 (df = 30)	2.9 (df = 29)
F Statistic	76.5*** (df = 1; 30)	91.4*** (df = 1; 30)	51.7*** (df = 2; 29)

Note:

*p<0.1; **p<0.05; ***p<0.01

Data are courtesy of Motor Trend Magazine. Also, R rules!