

Tutorial 4: Regression Model Estimation

Anh Le (anh.le@duke.edu)

September 18, 2015

Agenda

1. Create data frames
 2. Subset data frames
 3. Estimate a linear model with `lm()`
 4. Tips and tricks
- Prefix your objects in R (and related TAB tricks, i.e. arguments within function, variables within a data frame)
 - `fig.height()`, `fig.width()`
 - Code length ≤ 80 (RStudio > Tools > Options > Code)

1. Create data frames

```
my_dataframe <- data.frame(var1 = c(11, 12, 13),
                           var2 = c(21, 22, 23),
                           var3 = c("a", "b", "c"))
my_dataframe
```

```
##   var1 var2 var3
## 1   11   21    a
## 2   12   22    b
## 3   13   23    c
```

2. Subset data frames

All subsetting can be done with the following construct: `my_dataframe[?1 , ?2]`

The first question mark (?1) refers to which rows we want. The second question mark (?2) refers to which columns we want.

How to indicate to R which rows / columns we want? Multiple ways:

1. Use rows / columns index

```
my_dataframe[1, 2]
```

```
## [1] 21
```

```
my_dataframe[1:2, 2]
```

```
## [1] 21 22
```

```
my_dataframe[1:2, ]
```

```
##   var1 var2 var3
## 1   11   21    a
## 2   12   22    b
```

Rapid fire quiz

```
my_dataframe[2:3, ]
my_dataframe[ , 1:2]
my_dataframe[1:2, 2:3]

my_dataframe[c(1, 3), ]
my_dataframe[c(1, 3, 2), ]
```

2. Use rows / columns name

```
my_dataframe[ , "var2"]
```

```
## [1] 21 22 23
```

Rapid fire quiz:

```
my_dataframe[ , c("var1", "var3")]
my_dataframe[c(2, 3), c("var1", "var3")]
```

3. Use a condition

```
my_dataframe[c(TRUE, TRUE, FALSE), ]
```

```
##   var1 var2 var3
## 1   11   21    a
## 2   12   22    b
```

```
my_dataframe[, c(TRUE, FALSE, TRUE)]
```

```
##   var1 var3
## 1   11    a
## 2   12    b
## 3   13    c
```

Of course this is not tenable for a large data frame. So we have this very useful trick:

```
my_dataframe[my_dataframe$var1 < 13, ]
```

```
##   var1 var2 var3
## 1   11   21   a
## 2   12   22   b
```

This works because `my_dataframe$var1 < 13` actually returns `c(TRUE, TRUE, FALSE)` (vectorized operation in the wild!). Indeed:

```
my_dataframe$var1 < 13
```

```
## [1]  TRUE  TRUE FALSE
```

Rapid fire quiz:

```
my_dataframe[my_dataframe$var2 == 22, ]
my_dataframe[my_dataframe$var2 == 25, ]
```

4. Use a combination of condition

```
my_dataframe[my_dataframe$var1 > 10 & my_dataframe$var2 > 21, ]
```

```
##   var1 var2 var3
## 2   12   22   b
## 3   13   23   c
```

```
my_dataframe[my_dataframe$var1 > 10 | my_dataframe$var2 > 21, ]
```

```
##   var1 var2 var3
## 1   11   21   a
## 2   12   22   b
## 3   13   23   c
```

3. Estimate a linear model with `lm()`

In this section, I'll demo a (simplified) pipeline of steps in doing regression analysis with real data.

Download and clean data

```
library(WDI)
```

```
## Loading required package: RJSONIO
```

```
d_2010 <- WDI(indicator = c("NY.GDP.PCAP.CD", "SP.DYN.IMRT.IN", "SH.MED.PHYS.ZS"),
              start = 2010, end = 2010, extra = TRUE)
# d_2010[d_2010$region != "Aggregates", ]
```

There are a lot of unwanted columns. What if I just want `country`, `year`, and the three variables of interest (NY.GDP.PCAP.CD, SP.DYN.IMRT.IN, SH.MED.PHYS.ZS)? (Hint: subsetting)

```
d_2010 <- d_2010[d_2010$region != "Aggregates",
                 c("country", "year", "NY.GDP.PCAP.CD", "SP.DYN.IMRT.IN", "SH.MED.PHYS.ZS")]
```

Rename columns:

```
colnames(d_2010)
```

```
## [1] "country"      "year"          "NY.GDP.PCAP.CD" "SP.DYN.IMRT.IN"
## [5] "SH.MED.PHYS.ZS"
```

```
colnames(d_2010)[3:5] <- c('gdppc', 'infant_mortality', 'number_of_physician')
colnames(d_2010)
```

```
## [1] "country"      "year"          "gdppc"
## [4] "infant_mortality" "number_of_physician"
```

Log gdp per capita

```
d_2010$log_gdppc <- log(d_2010$gdppc)
```

Build a linear model

```
lm(infant_mortality ~ log_gdppc, data = d_2010)
```

```
##
## Call:
## lm(formula = infant_mortality ~ log_gdppc, data = d_2010)
##
## Coefficients:
## (Intercept)    log_gdppc
##      139.31         -13.14
```

```
m1 <- lm(infant_mortality ~ log_gdppc, data = d_2010)
summary(m1)
```

```
##
## Call:
## lm(formula = infant_mortality ~ log_gdppc, data = d_2010)
##
## Residuals:
```

```
##      Min      1Q  Median      3Q      Max
## -27.641 -10.393  -2.596   6.222  79.971
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 139.3147     6.8043   20.47  <2e-16 ***
## log_gdppc   -13.1403     0.7909  -16.61  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 16.14 on 184 degrees of freedom
## (31 observations deleted due to missingness)
## Multiple R-squared:  0.6, Adjusted R-squared:  0.5979
## F-statistic: 276 on 1 and 184 DF, p-value: < 2.2e-16
```

```
m2 <- lm(infant_mortality ~ log_gdppc + number_of_physician, data = d_2010)
summary(m2)
```

```
##
## Call:
## lm(formula = infant_mortality ~ log_gdppc + number_of_physician,
##     data = d_2010)
##
## Residuals:
##      Min      1Q  Median      3Q      Max
## -27.343  -8.182  -1.260   5.961  50.581
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    121.1074     7.4312  16.297  < 2e-16 ***
## log_gdppc      -10.5850     0.9824 -10.774  < 2e-16 ***
## number_of_physician -2.9102     0.9288  -3.133  0.00211 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.71 on 140 degrees of freedom
## (74 observations deleted due to missingness)
## Multiple R-squared:  0.7078, Adjusted R-squared:  0.7036
## F-statistic: 169.6 on 2 and 140 DF, p-value: < 2.2e-16
```

Extract result from the model

`str()` (stands for structure) is used to look into the structure of an object in R, see what it contains.

```
str(m1)
```

```
## List of 13
## $ coefficients : Named num [1:2] 139.3 -13.1
##   .. attr(*, "names")= chr [1:2] "(Intercept)" "log_gdppc"
## $ residuals    : Named num [1:186] 3.36 5.05 18.96 -7.12 -15.22 ...
##   .. attr(*, "names")= chr [1:186] "6" "7" "8" "9" ...
## $ effects      : Named num [1:186] -381.68 268.13 19.23 -7.54 -15.39 ...
```

```

## ..- attr(*, "names")= chr [1:186] "(Intercept)" "log_gdppc" "" "" ...
## $ rank : int 2
## $ fitted.values: Named num [1:186] -0.864 2.251 56.135 14.823 30.022 ...
## ..- attr(*, "names")= chr [1:186] "6" "7" "8" "9" ...
## $ assign : int [1:2] 0 1
## $ qr :List of 5
## ..$ qr : num [1:186, 1:2] -13.6382 0.0733 0.0733 0.0733 0.0733 ...
## .. ..- attr(*, "dimnames")=List of 2
## .. ..$ : chr [1:186] "6" "7" "8" "9" ...
## .. ..$ : chr [1:2] "(Intercept)" "log_gdppc"
## .. ..- attr(*, "assign")= int [1:2] 0 1
## ..$ qraux: num [1:2] 1.07 1.09
## ..$ pivot: int [1:2] 1 2
## ..$ tol : num 1e-07
## ..$ rank : int 2
## ..- attr(*, "class")= chr "qr"
## $ df.residual : int 184
## $ na.action :Class 'omit' Named int [1:31] 9 12 14 24 45 46 61 65 72 78 ...
## .. ..- attr(*, "names")= chr [1:31] "14" "17" "NA" "29" ...
## $ xlevels : Named list()
## $ call : language lm(formula = infant_mortality ~ log_gdppc, data = d_2010)
## $ terms :Classes 'terms', 'formula' length 3 infant_mortality ~ log_gdppc
## .. ..- attr(*, "variables")= language list(infant_mortality, log_gdppc)
## .. ..- attr(*, "factors")= int [1:2, 1] 0 1
## .. .. ..- attr(*, "dimnames")=List of 2
## .. .. ..$ : chr [1:2] "infant_mortality" "log_gdppc"
## .. .. ..$ : chr "log_gdppc"
## .. ..- attr(*, "term.labels")= chr "log_gdppc"
## .. ..- attr(*, "order")= int 1
## .. ..- attr(*, "intercept")= int 1
## .. ..- attr(*, "response")= int 1
## .. ..- attr(*, ".Environment")=<environment: R_GlobalEnv>
## .. ..- attr(*, "predvars")= language list(infant_mortality, log_gdppc)
## .. ..- attr(*, "dataClasses")= Named chr [1:2] "numeric" "numeric"
## .. .. ..- attr(*, "names")= chr [1:2] "infant_mortality" "log_gdppc"
## $ model :'data.frame': 186 obs. of 2 variables:
## ..$ infant_mortality: num [1:186] 2.5 7.3 75.1 7.7 14.8 ...
## ..$ log_gdppc : num [1:186] 10.67 10.43 6.33 9.47 8.32 ...
## ..- attr(*, "terms")=Classes 'terms', 'formula' length 3 infant_mortality ~ log_gdppc
## .. .. ..- attr(*, "variables")= language list(infant_mortality, log_gdppc)
## .. .. ..- attr(*, "factors")= int [1:2, 1] 0 1
## .. .. .. ..- attr(*, "dimnames")=List of 2
## .. .. .. ..$ : chr [1:2] "infant_mortality" "log_gdppc"
## .. .. .. ..$ : chr "log_gdppc"
## .. .. ..- attr(*, "term.labels")= chr "log_gdppc"
## .. .. ..- attr(*, "order")= int 1
## .. .. ..- attr(*, "intercept")= int 1
## .. .. ..- attr(*, "response")= int 1
## .. .. ..- attr(*, ".Environment")=<environment: R_GlobalEnv>
## .. .. ..- attr(*, "predvars")= language list(infant_mortality, log_gdppc)
## .. .. ..- attr(*, "dataClasses")= Named chr [1:2] "numeric" "numeric"
## .. .. .. ..- attr(*, "names")= chr [1:2] "infant_mortality" "log_gdppc"
## ..- attr(*, "na.action")=Class 'omit' Named int [1:31] 9 12 14 24 45 46 61 65 72 78 ...
## .. .. ..- attr(*, "names")= chr [1:31] "14" "17" "NA" "29" ...

```

```
## - attr(*, "class")= chr "lm"
```

```
m1$coefficients
```

```
## (Intercept)  log_gdppc  
##   139.31467   -13.14033
```

```
m1$coefficients['(Intercept)']
```

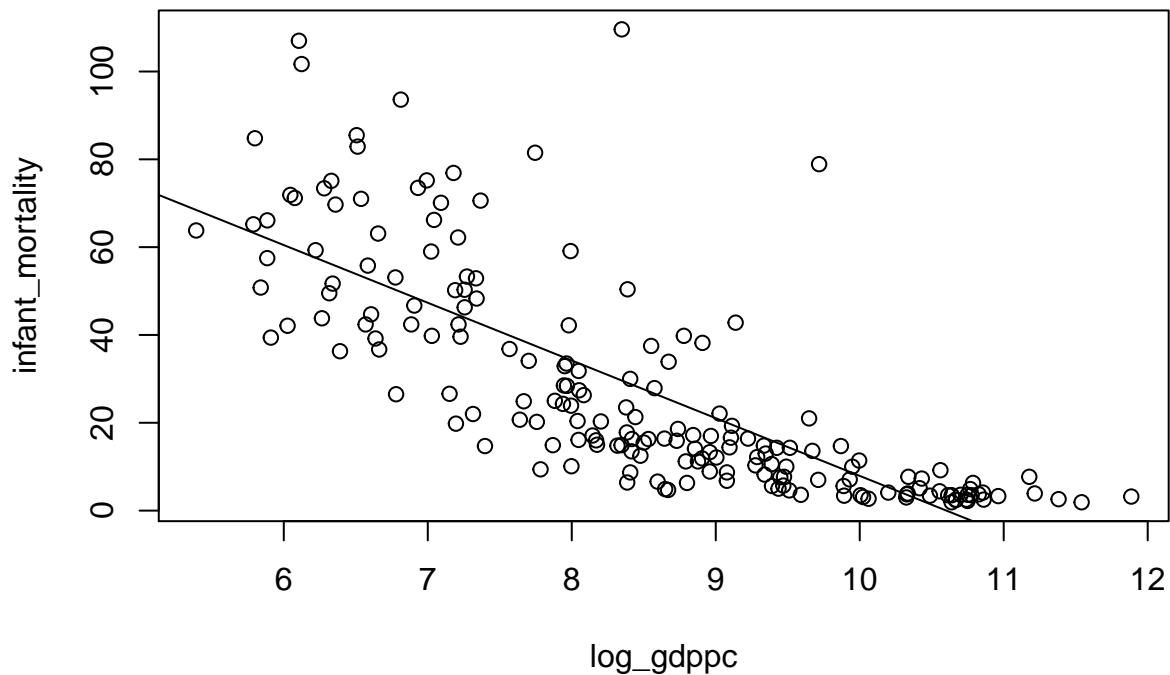
```
## (Intercept)  
##   139.3147
```

```
m1$coefficients['log_gdppc']
```

```
## log_gdppc  
## -13.14033
```

Now we can use them for other things, e.g plotting

```
plot(infant_mortality ~ log_gdppc, data = d_2010)  
abline(a = m1$coefficients['(Intercept)'], b = m1$coefficients['log_gdppc'])
```



Report the model in a nice, journal-ready format

The **stargazer** library takes your model objects and generates tables in LaTeX. This package has a lot of customizing options, which you'll explore in the homework.

```
library(stargazer)
```

```
##  
## Please cite as:  
##  
## Hlavac, Marek (2014). stargazer: LaTeX code and ASCII text for well-formatted regression and summary  
## R package version 5.1. http://CRAN.R-project.org/package=stargazer
```

```
# LaTeX code that you can copy paste into LaTeX  
stargazer(m1, m2)
```

```
##  
## % Table created by stargazer v.5.1 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu  
## % Date and time: Tue, Sep 15, 2015 - 01:13:26 PM  
## \begin{table}[!htbp] \centering  
## \caption{}  
## \label{}  
## \begin{tabular}{@{\extracolsep{5pt}}lcc}  
## \hline  
## \hline \hline  
## & \multicolumn{2}{c}{\textit{Dependent variable:}} \\\br/>## \cline{2-3}  
## \hline & \multicolumn{2}{c}{infant\_mortality} \\\br/>## \hline & (1) & (2) \\\br/>## \hline  
## log\_gdppc &  $-13.140^{***}$  &  $-10.585^{***}$  \\\br/>## & (0.791) & (0.982) \\\br/>## & & \\\br/>## number\_of\_physician &  $-2.910^{***}$  & \\\br/>## & (0.929) \\\br/>## & & \\\br/>## Constant &  $139.315^{***}$  &  $121.107^{***}$  \\\br/>## & (6.804) & (7.431) \\\br/>## & & \\\br/>## \hline \hline  
## Observations & 186 & 143 \\\br/>## R2 & 0.600 & 0.708 \\\br/>## Adjusted R2 & 0.598 & 0.704 \\\br/>## Residual Std. Error & 16.138 (df = 184) & 12.711 (df = 140) \\\br/>## F Statistic & 276.047*** (df = 1; 184) & 169.553*** (df = 2; 140) \\\br/>## \hline  
## \hline \hline  
## \textit{Note:} & \multicolumn{2}{r}{ $^{*}p < 0.1$ ;  $^{**}p < 0.05$ ;  $^{***}p < 0.01$ } \\\br/>## \end{tabular}  
## \end{table}
```

```
# If using knitr, use the option results='asis'  
stargazer(m1, m2)
```

```
% Table created by stargazer v.5.1 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu  
% Date and time: Tue, Sep 15, 2015 - 01:13:26 PM
```


Table 1:

	<i>Dependent variable:</i>	
	infant_mortality	
	(1)	(2)
log_gdppc	-13.140*** (0.791)	-10.585*** (0.982)
number_of_physician		-2.910*** (0.929)
Constant	139.315*** (6.804)	121.107*** (7.431)
Observations	186	143
R ²	0.600	0.708
Adjusted R ²	0.598	0.704
Residual Std. Error	16.138 (df = 184)	12.711 (df = 140)
F Statistic	276.047*** (df = 1; 184)	169.553*** (df = 2; 140)
<i>Note:</i>		

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