

Lab 3: Regression Model Estimation

Anh Le (*anh.le@duke.edu*)

September 21, 2018

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)

## Warning: package 'WDI' was built under R version 3.4.4
## 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)
```

Remove missing data

```
d_2010 <- na.omit(d_2010)
```

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
##      134.29      -12.78

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
## -24.938  -9.380  -1.190   7.333  50.665
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   134.294     6.519   20.60  <2e-16 ***
## log_gdppc     -12.783     0.757  -16.89  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 13.39 on 146 degrees of freedom
## Multiple R-squared:  0.6614, Adjusted R-squared:  0.659
## F-statistic: 285.1 on 1 and 146 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
## -24.935  -7.842  -1.379   6.957  51.240
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   116.760     7.590  15.384  < 2e-16 ***
## log_gdppc     -9.935     1.011  -9.823  < 2e-16 ***
## number_of_physician -4.047     1.009  -4.013  9.58e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.75 on 145 degrees of freedom
```

```
## Multiple R-squared:  0.6952, Adjusted R-squared:  0.691
## F-statistic: 165.4 on 2 and 145 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 12
## $ coefficients : Named num [1:2] 134.3 -12.8
##   ..- attr(*, "names")= chr [1:2] "(Intercept)" "log_gdppc"
## $ residuals    : Named num [1:148] 4.87 6.67 12.04 -17.28 -14.95 ...
##   ..- attr(*, "names")= chr [1:148] "6" "7" "8" "10" ...
## $ effects      : Named num [1:148] -313.7 226.1 12.4 -17.5 -15.1 ...
##   ..- attr(*, "names")= chr [1:148] "(Intercept)" "log_gdppc" "" "" ...
## $ rank         : int 2
## $ fitted.values: Named num [1:148] -1.074 0.535 53.56 27.977 31.054 ...
##   ..- attr(*, "names")= chr [1:148] "6" "7" "8" "10" ...
## $ assign       : int [1:2] 0 1
## $ qr           :List of 5
##   ..$ qr      : num [1:148, 1:2] -12.1655 0.0822 0.0822 0.0822 0.0822 ...
##   .. ..- attr(*, "dimnames")=List of 2
##   .. .. ..$ : chr [1:148] "6" "7" "8" "10" ...
##   .. .. ..$ : chr [1:2] "(Intercept)" "log_gdppc"
##   .. ..- attr(*, "assign")= int [1:2] 0 1
##   ..$ qraux: num [1:2] 1.08 1.1
##   ..$ pivot: int [1:2] 1 2
##   ..$ tol   : num 1e-07
##   ..$ rank  : int 2
##   ..- attr(*, "class")= chr "qr"
## $ df.residual  : int 146
## $ xlevels      : Named list()
## $ call         : language lm(formula = infant_mortality ~ log_gdppc, data = d_2010)
## $ terms        :Classes 'terms', 'formula' language 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':  148 obs. of  2 variables:
##   ..$ infant_mortality: num [1:148] 3.8 7.2 65.6 10.7 16.1 13 3.6 32.6 6.2 13.4 ...
##   ..$ log_gdppc       : num [1:148] 10.59 10.46 6.32 8.32 8.08 ...
##   ..- attr(*, "terms")=Classes 'terms', 'formula' language infant_mortality ~ log_gdppc
##   .. .. ..- attr(*, "variables")= language list(infant_mortality, log_gdppc)
##   .. .. ..- attr(*, "factors")= int [1:2, 1] 0 1
```

You can extract the coefficients

```
## (Intercept)    log_gdppc
##    134.29374    -12.78254
```

```
## (Intercept)
##      134.2937
```

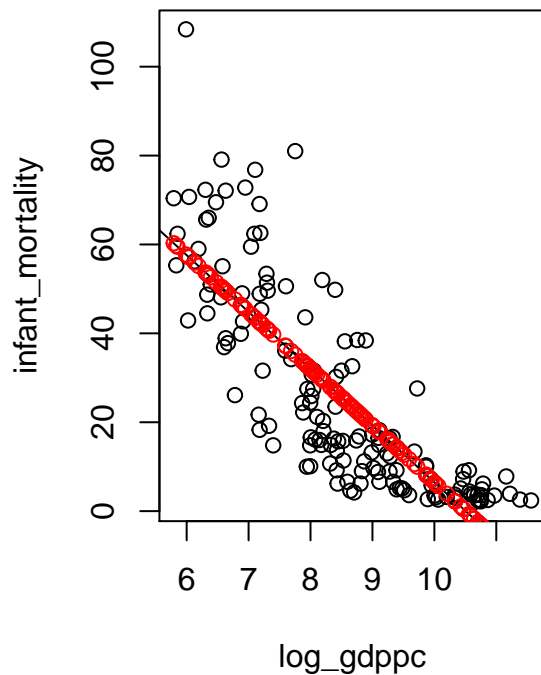
```
## log_gdppc
## -12.78254
```

```
d_2010$pred_infant_mortality2 <- m1$coefficients['(Intercept)'] + m1$coefficients['log_gdppc'] * d_2010
```

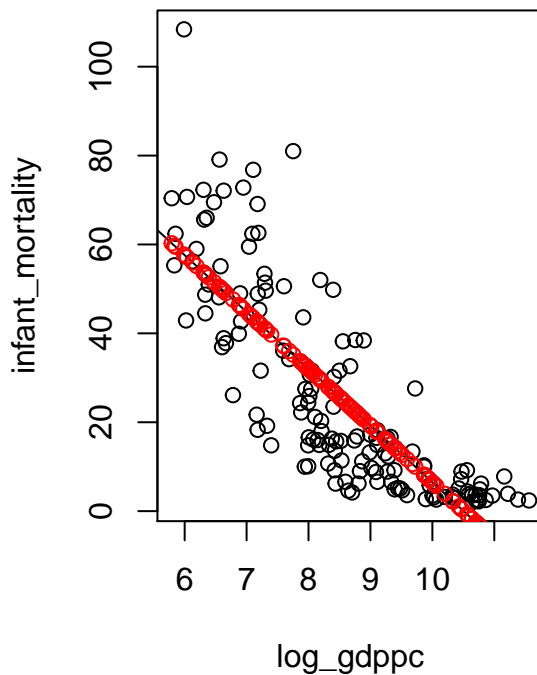
```
par(mfrow = c(1, 2))
plot(infant_mortality ~ log_gdppc, data = d_2010, main = "Plot predicted values-method 1")
abline(a = m1$coefficients['(Intercept)'], b = m1$coefficients['log_gdppc'])
points(d_2010$log_gdppc, d_2010$pred_infant_mortality1, col = 'red')

plot(infant_mortality ~ log_gdppc, data = d_2010, main = "Plot predicted values-method 2")
abline(a = m1$coefficients['(Intercept)'], b = m1$coefficients['log_gdppc'])
points(d_2010$log_gdppc, d_2010$pred_infant_mortality2, col = 'red')
```

Plot predicted values–method 1



Plot predicted values–method 2



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)
```

```
## Warning: package 'stargazer' was built under R version 3.4.4
```

```
##
```

```
## Please cite as:
```

```
## Hlavac, Marek (2018). stargazer: Well-Formatted Regression and Summary Statistics Tables.
```

```
## R package version 5.2.2. https://CRAN.R-project.org/package=stargazer
```

```
# LaTeX code that you can copy paste into LaTeX
```

```
stargazer(m1, m2)
```

```
##
```

```
## % 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 21, 2018 - 12:11:49 PM
```

```
## \begin{table}[!htbp] \centering
```

```
## \caption{}
```

```
## \label{}
```

```
## \begin{tabular}{@{\extracolsep{5pt}}lcc}
```

```
## \hline
```

```
## \hline \hline
```

```

## & \multicolumn{2}{c}{\textit{Dependent variable:}} \\
## \cline{2-3}
## \[-1.8ex] & \multicolumn{2}{c}{infant\_mortality} \\
## \[-1.8ex] & (1) & (2) \\
## \hline \[-1.8ex]
## log\_gdppc & $-12.783^{***}$ & $-9.935^{***}$ \\
## & (0.757) & (1.011) \\
## & & \\
## number\_of\_physician & & $-4.047^{***}$ \\
## & & (1.009) \\
## & & \\
## Constant & 134.294^{***} & 116.760^{***} \\
## & (6.519) & (7.590) \\
## & & \\
## \hline \[-1.8ex]
## Observations & 148 & 148 \\
## R^2 & 0.661 & 0.695 \\
## Adjusted R^2 & 0.659 & 0.691 \\
## Residual Std. Error & 13.391 (df = 146) & 12.748 (df = 145) \\
## F Statistic & 285.146^{***} (df = 1; 146) & 165.372^{***} (df = 2; 145) \\
## \hline
## \hline \[-1.8ex]
## \textit{Note:} & \multicolumn{2}{c}{\textit{$^*p<0.1$; $^{**}p<0.05$; $^{***}p<0.01$}} \\
## \end{tabular}
## \end{table}

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

```

% 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 21, 2018 - 12:11:54 PM

Table 1:		
	<i>Dependent variable:</i>	
	infant_mortality	
	(1)	(2)
log_gdppc	-12.783*** (0.757)	-9.935*** (1.011)
number_of_physician		-4.047*** (1.009)
Constant	134.294*** (6.519)	116.760*** (7.590)
Observations	148	148
R ²	0.661	0.695
Adjusted R ²	0.659	0.691
Residual Std. Error	13.391 (df = 146)	12.748 (df = 145)
F Statistic	285.146*** (df = 1; 146)	165.372*** (df = 2; 145)
<i>Note:</i>		
*p<0.1; **p<0.05; ***p<0.01		