

Hands-On Data Analysis for ININ Using R

Prof. Dr. Cornelia Storz, M.Sc. Fei (Michael) Wang

Management and Microeconomics, Goethe-Universität Frankfurt

This document was prepared for students who are taking ININ course and planning to take the exam. It is a collection of notes and codes for the course. The notes are based on the tutorials we had in the course. I am trying to make it concise and easy to understand. I hope it can help you to review the course and prepare for the exam. *We are living in a very noisy world, therefore let's keep it simple and clear.* I setup a challenge for myself to deliver a clear and concise review notes within 15 pages. This brings the trade-off, which means some figures and tables are not included in the notes. Therefore, you have to run the codes to see the results.

I hope you enjoy reading it. I also hope you will have this notes with you whenever you want to do some data analysis. If one day, you still refer to this notes and find it still useful, I would be very happy to hear that.

Keywords: econometrics, data analysis, regression models, empirical research, innovation, management

Contents

1 Introduction	1
2 Introduction to Data and <code>data.table</code>	1
2.1 Data Visualization	3
2.2 Log Transformation	4

1 Introduction

All statistical or econometric or machine learning models are based on the following assumptions:


- there are something we know - **data**
- and something we don't know - **error ϵ** .

In summary, according to confucius, *to know what we know and what we do not know*, that is called **wisdom**. Or like Plato said, *I know that I know nothing*. To help you to review the course, the notes will be organized as follows:

1. **Data**: using `data.table` to get familiar with the data
2. **Simple linear regression**: how to estimate a simple linear regression model, how to interpret the results
3. **Multiple linear regression**: how to estimate a multiple linear regression model, how to interpret the results, how to test the model
4. **Introduction to logistic regression**: why do we need logistic regression
5. **Data manipulation**: will not be tested in the exam, but it is very useful for your future work or research

2 Introduction to Data and `data.table`

Broadly speaking, there are two kinds of data: **structured data** and **unstructured data**. Structured data is data that has a structure, such as a table, whereas unstructured data is data that does not have a structure, such as a text file. In this course, we focus on structured data. This means all the data we will use look like tables, such as the following one:



oceannumeric.github.io

```
# create a data.table
dt <- data.table(
  vn1 = c(1, 100, -567),
  vn2 = c("hello", "hello", "hello"),
  vn3 = rep("world", 3)
)
```

```
# read a csv file into data.table
dt <- fread("file_name.csv")
```

`dt[i, j, by]`

any operation on columns takes place at *j*

variable_name_1	variable_name_2	variable_name_3	variable_name_4	vn5	vn7	vn8
integer	numeric (dbl)	character	factor	logic	mixed with missing values	Date/Time
1	2.0	A	female / 1	TRUE	2.0	2017-09-16
100	-3.1415926	"hello"	male / 2	FALSE	"abc"	16:23:57
-567	100	hello world	any categorical data	TRUE	NA	2 June 2020

any operation on rows takes place at *i*

common functions:

```
str(dt)
summary(dt)
names(dt)
dim(dt)
```

```
#save a data.table into a csv file
fwrite(dt, "file_name.csv")
```

The basic syntax of `data.table` is summarized in the following illustration. **You will not be tested on the syntax of `data.table` in the exam.** However, you will be tested on the underlying concepts of

data.table, such as the type of variables (integer, character, factor, etc.). In the future if you will be working as a data scientist, you can use data.table to do big data analysis. You will need to know the syntax of data.table for practical use not for the exam.

Import key packages

```
library(data.table)
library(magrittr)
library(knitr)
library(ggplot2)
```

Structure of the dataset

```
str(dt)
head(dt)
summary(dt)
names(dt)
setnames('old', 'new')
setorder(vn5, vn6)
apply(dt, function(x) sum(is.na(x)))
```

Check unique or duplicated values

```
dt %>%
  unique(by = c("variables"))
dt %>%
  .[duplicated(variable)]
# print out all duplicates
dt %>%
  .[duplicated(variable) | duplicated(variable, fromLast = TRUE)]
```



dt[i, j, by]

```
# common functions in pipe
%>%
with()
kable()
plot()
par(mfrow = c(2, 2))
# ask ChatGPT always
```

dt						
variable_name_1	variable_name_2	variable_name_3	variable_name_4	vn5	vn7	vn8
integer	numeric (dbl)	character	factor	logic	mixed with missing	Date/Time
1	2.0	A	female / 1	TRUE	2.0	2013
100	-3.1415926	hello	male / 2	FALSE	"abc"	2022-09-10
-567	100	hello world	any categorical data	TRUE	NA	09:12:37

```
# class
class(variable)
# is function
is.factor()
is.integer()
is.character()
# as function
as.character()
as.integer()
as.POSIXct()
as.factor()
...
```

Manipulate rows with i

```
# extract rows based on index
dt[5:17, ] # all columns
dt[1:9, 2:4] # row 1 to 9, column 2 to 4
# subset rows based on conditions
dt %>%
  .[vn2 >= 20]

# logical operators to use in i
> < >= <= is.na() !
& | %in% %like% %between%

# any functions from dplyr
# could be combined within
# the pipe line
```

Manipulate columns with j

```
# extract columns
dt %>%
  .[, .(vn5, vn7)]
dt %>%
  .[, c(2:6)] # using column index
# extract columns based on names
dt %>%
  .[, .SD, .SDcols = patterns("^a")]
# extract columns based on type
dt %>%
  .[, .SD, .SDcols = is.integer]
# extract and transform at the same time
dt %>%
  .[, lapply(.SD, tolower), .SDcols = is.character]
# create a new columns on original data
dt %>%
  .[, vn9 := vn2 + 2]
# create or transform columns on original data
# using name vector, .SDcols, and lapply with :=
```

Subgroup with by

```
# summarize vn2 by vn4
dt %>%
  .[, .(vn2_mean = mean(vn2), by = vn4)]

# one of the most common way to use by
# is that we need do some operation on one
# or several variables based on another
# categorical variables, such as
- dt[, .(c=sum(b)), by = a]
- dt[, .(c=mean(b)), by = .(a, d)]

# use keyby if you want to
# sort within the group

# special functions that
# could bring magics
.N, .I .SD
```

Let's start from installing and loading the packages we need for the course.

```
1 # install packages
2 install.packages("stargazer")
3 # install ISLR if you don't have it
4 # install.packages("ISLR")
5 install.packages("corrplot")
6 # sometimes you need to install other packages
7 # hopefully you can figure it out by yourself
```

After you install the packages, you can load them into R.

```
1 # library for data analysis
2 library(data.table)
3 library(magrittr)
4 library(ggplot2)
5 library(knitr)
6 library(stargazer)
7 library(MASS)
8 library(ISLR)
9 library(corrplot)
```

Now, we can load the data into R and manipulate the data.

```
1 # read the dataset
2 cis <- fread("https://shorturl.at/wBESZ")
3 # check structure of the dataset
4 str(cis)
5 # check the first 10 rows of the dataset
6 head(cis, 10)
```

```

7 # check the number of missing values in each column
8 cis %>%
9   .[, .SD, .SDcols = is.double] %>%
10   # check the number of missing values in each column
11   apply(function(x) sum(is.na(x))) %>%
12   # sort the number of missing values in each column
13   sort(decreasing = TRUE) %>%
14   # convert to data.table and keep rownames as a column
15   as.data.table(keep.rownames = TRUE) %>%
16   # set variable names
17   setnames(c("variables", "Numbe of NAs")) %>%
18   # check the first 10
19   head(10) %>%
20   kable("pipe")

```

2.1 Data Visualization

It is important to visualize the data before you start to do the analysis. To choose the right figure, you need to know the type of variables. Here is the summary:

- **Categorical variables:** bar chart, pie chart, etc.
- **Continuous variables:** histogram, boxplot, etc.
- **Categorical vs. continuous variables:** boxplot or violin or histogram with different colors

Here is a demo of how to visualize the data. You can try to run the code and see the results.

```

1 # four figures in one page
2 # bar chart, histogram, box plot and box plot compared
3 # figure size
4 options(repr.plot.width = 10, repr.plot.height = 10)
5 par(mfrow = c(2, 2))
6 cis %>%
7   # group by branche and .N calculates the number of frequency
8   .[, .N, by = branche] %>%
9   # order it in a descending way
10  .[order(-N)] %>%
11  # get the top 10 branche (industry)
12  head(10) %>%
13  # plot it
14  with(pie(N, labels = branche, main = "Distribution of Industries"))
15
16 # histgorma for number of employees
17 cis %>%
18   with(hist(bges, main = "number of employees"))
19 # boxplot for sales
20 cis %>%
21   with(boxplot(um18, main = "Boxplot of Sales in 2018"))
22 # boxplot for log(1+sales)
23 cis %>%
24   with(boxplot(log(1+um18), main = "Boxplot of Log Transform "))

```

When you visualize the data, it is important to check two things for continuous variables:

- **shape:** whether the distribution is symmetric or skewed (ideally, we prefer symmetric distribution)
- **outliers:** outliers are extreme values that deviate from other observations on data , they may indicate a variability in a measurement, experimental errors or a novelty.

2.2 Log Transformation

Log transformation is a very useful tool to deal with skewed data. When you have a skewed data, you can try to log transform it. However, it could be tricky. You need to know the underlying theory of log transformation. Generally speaking, log transformation is used to make the data more symmetric. To avoid the negative values, we usually use $\log(1+x)$ instead of $\log(x)$.

```
1 # log transform
2 options(repr.plot.width = 10, repr.plot.height = 10)
3 par(mfrow = c(2, 2))
4 cis %>%
5   with(hist(um18, main = "Histogram of Sales (2018)"))
6
7 cis %>%
8   with(hist(log(1+um18), main = "Histogram of Log (1+sales)"))
9
10 cis %>%
11   with(boxplot(um18, main = "Boxplot of Sales (2018)"))
12 cis %>%
13   with(boxplot(log(1+um18), main = "Boxplot of Log Transform "))
```

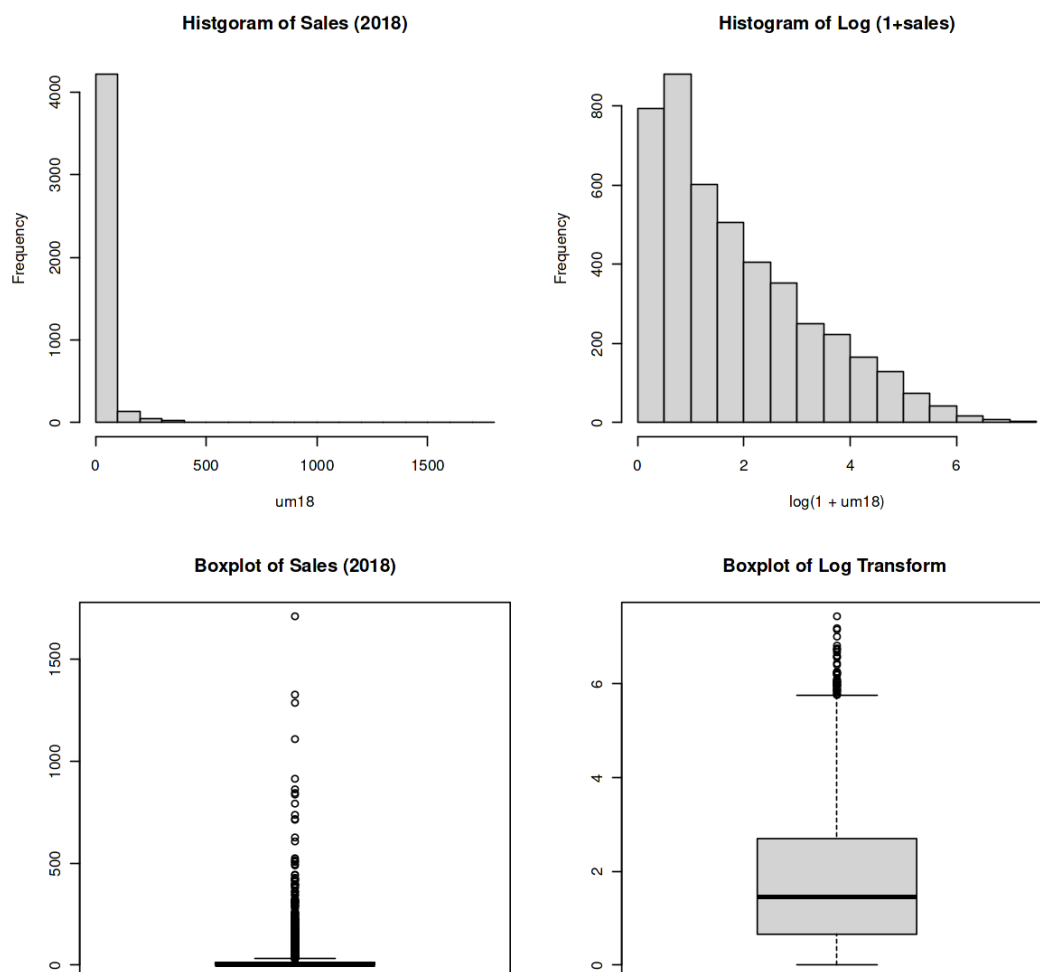


Figure 1: Log transformation of sales