```
1    # Importing Data
2
3    # The first step for data management and alaysis is to have data. You may
4    # enter data using the keyboard, but, in most cases, you will import data file
5    # (from Excel, Text, Stata, Sas, etc) into R.
6
7    # 1. csv: comman separated values
8    getwd()
9    setwd("C:/Users/Min Seong Kim/Dropbox/R_programming/lecture/elsect_main")
10   dir()
11   rev_exp0 <- read.csv("district_rev_exp.csv")
12
13   rev_exp <- read.csv(file.choose(), stringsAsFactors = FALSE)  # district_rev_exp.csv
14
15   rev_exp$STATE <- as.factor(rev_exp$STATE)
16
17   class(rev_exp)
18   head(rev_exp)
19   tail(rev_exp)
20
21   str(rev_exp)
22   summary(rev_exp)
23
24   rev_exp$TOTALREV[rev_exp$TOTALREV == "-"] <- NA
25   rev_exp$TOTALREV <- as.numeric(rev_exp$TOTALREV)
26
27   # Calculate the averages of enrollment, total revenue and total expenditure for each
     state.
28   colMeans(rev_exp[rev_exp$STATE == "California",c("ENROLL", "TOTALREV", "TOTALEXP")],
     na.rm=TRUE)
29   colMeans(rev_exp[rev_exp$STATE == "Connecticut",c("ENROLL", "TOTALREV", "TOTALEXP")])
30   colMeans(rev_exp[rev_exp$STATE == "Massachusetts",c("ENROLL", "TOTALREV", "TOTALEXP")])
31   colMeans(rev_exp[rev_exp$STATE == "Missouri",c("ENROLL", "TOTALREV", "TOTALEXP")])
32
33   aggregate(rev_exp[ ,c("ENROLL", "TOTALREV", "TOTALEXP")], list(rev_exp$STATE), mean,
     na.rm=TRUE)
34   # list() specifies the criterion to make groups
35   aggregate(rev_exp[ ,c("ENROLL", "TOTALREV", "TOTALEXP")],
36           by=list(ST = rev_exp$STATE, EnR = rev_exp$ENROLL > 1000), FUN=mean,
               na.rm=TRUE)
37
38   a <- na.omit(rev_exp)  # eleminate the rows that contains NA.
39
40   # 2. text file:  tab-delimited file
41   rev_exp1 <- read.delim(file.choose(), stringsAsFactors = FALSE)   # district_rev_exp.txt
42
43   head(rev_exp1)
44   tail(rev_exp1)
45
46   # 3. read.table: read any tabular file as a data.frame
47
48   # Use district_rev_exp_readtable.txt
49   rev_exp2 <- read.table(file.choose(), sep="/", stringsAsFactors = FALSE)
50   names(rev_exp2)
51   rev_exp2 <- read.table(file.choose(), sep="/", header = TRUE, stringsAsFactors = FALSE)
52   names(rev_exp2)
53
54   # We can also read csv file and tab delimited txt file using read.table
55   rev_exp3 <- read.table(file.choose(), sep=",", header = TRUE, stringsAsFactors =
     FALSE)   # district_rev_exp.csv
56   rev_exp3 <- read.table(file.choose(), sep="\t", header = TRUE, stringsAsFactors =
     FALSE)   # district_rev_exp.csv
57
58   # You can save excel file with csv or tab delimited txt file. Then, you can use the
     functions above
59   # to read the file.
60   # You can read excel file directly.
61   # First install the package "readxl"
62   install.packages(readxl)
```

```r
63   library(readxl)
64   excel_sheets(file.choose())                          # list different sheets
65   rev_exp4 <- read_excel(file.choose(), sheet = 1)     # actually import data into R
66
67   # You can also import data from Stata
68   # You first install the package "foreign".
69   install.packages(foreign)
70   library(foreign)
71   read.dta(file.choose())
72
73   # Exercise
     ##############################################################################
74
75
76   # Download complete.csv from HustkyCT (lecture10-data) in your computer.
77   # Import this dataset to R.
78
79   # 1. Which league is the best in terms of wage (eur_wage) and overall?
80   # 2. Based on "eur_value", which team has the most players in top 100?
81   #    hint: 1. Sort based on eur_value, 2. Make sure that team is a factor
82   # 3. Present the distribution of average wage (eur_wage) based on age.
83   #    hint: 1. Use aggregate() function to obtain average wages for each age, 2. Use
     plot() function
84   # 4. Which team has the most players under 23?
85   #    hint: Make sure that team is a factor
86
87
88   sc <- read.csv(file.choose())
89   mean(sc$eur_wage)
90
```