

```

1  # Data structure: 4. Data frame
2  # A data frame is a two dimensional format of data structure that is useful to store
3  # data in tables. In general, each column is for a different item and
4  # each row is for a different unit.
5  # So, if you are interested in a specific item (for example, name or grade), you need to
6  # choose a column. If you are interested in information for a specific unit,
7  # you can choose a row.
8  # It is a collection of variables of the same length,
9  # but possibly of different type of variables.(numeric, factor, character)
10 # <-> Elements of a matrix should be of the same type.
11
12 rm(list=ls())
13 name <- c("Tom", "James", "Mary", "Paker")
14 score <- c(9, 7, 6, 10)
15 grade <- factor(c("A", "C", "D", "A"), ordered=TRUE, levels=c("D","C","B","A"))
16 Econ_dep <- data.frame(name, score, grade)
17
18 Econ_dep
19
20 str(Econ_dep) # gives us the number of observations (# of rows),
21               # the number of variables (# of columns), full list of the variable names,
22
23 # Name is not a factor (categorical variable). Two persons with the same name are not
24 # in the same category.
25 # A character vector is converted to a factor in data.frame if we don't specify as a
26 # character.
27 Econ_dep1 <- data.frame(name, score, grade, stringsAsFactors = FALSE)
28 str(Econ_dep1)
29 Econ_dep2 <- data.frame(I(name), score, grade)
30 str(Econ_dep2)
31
32 Econ_dep$name
33 name
34
35 rm(name, score, grade)
36
37 # subset selection and subset elimination
38 Econ_dep[c(3,1),]
39 Econ_dep[c(1,3),]
40 Econ_dep[-c(2,3),]
41 Econ_dep[c(TRUE,TRUE,FALSE,FALSE),]
42 Econ_dep[Econ_dep$name=="James",]
43 Econ_dep[!(Econ_dep$name=="James"),]
44 subset(x=Econ_dep, subset=!(Econ_dep$name=="James")) # subset()
45
46 Econ_dep[,2]
47 Econ_dep[,-2]
48 Econ_dep[,c("name","grade")]
49 Econ_dep[,c(TRUE, FALSE, TRUE)]
50 Econ_dep[,!c(TRUE, FALSE, TRUE)]
51
52 Econ_dep$grade <- NULL # eliminate "grade" column
53 Econ_dep
54
55 # Add a vector to data.frame
56 attendance <- c("all", "some", "never", "all")
57 Econ_dep <- cbind(Econ_dep,attendance)
58 Econ_dep$gender <- c("Male", "Male", "Female", "Female")
59 str(Econ_dep) # Econ_dep$gender is a character vector not a
60 # factor.
61
62 Econ_dep$year <- NA
63 Econ_dep$year[Econ_dep$name %in% c("Tom", "Mary")] <- 4
64 Econ_dep
65
66 # Sort midterm using order() function
67 rank <- order(Econ_dep$score, decreasing = TRUE)
68 Econ_dep[rank,]

```

```

67 attach(Econ_dep)           # data.frame "midterm" is attached to the R search path.
68                             # So, we don't need to put midterm$ to call a variable
69                             # in midterm
70
71 name
72
73 detach(Econ_dep)
74 name
75
76 as.matrix(Econ_dep)         # Easy to convert data.frame to matrix
77
78 # Grouped data
79 ID <- 1:20
80 rand.number <- rnorm(20)
81 participant <- data.frame(ID, rand.number)
82
83 # Generate two data frames "group1" and "group2" from "participant"
84
85 group1 <- participant[participant$rand.number > 0,]   # Conditional selection
86 group2 <- subset(participant, subset = rand.number <= 0) # You can use subset function.
87
88 # Create a new variable "level" in "participant" to group data.
89
90 participant$level[participant$rand.number > 0.5] <- "first"
91 participant$level[participant$rand.number > 0 & participant$rand.number <= 0.5] <-
  "second"
92 participant$level[participant$rand.number > -0.5 & participant$rand.number <= 0] <-
  "third"
93 participant$level[participant$rand.number <= -0.5] <- "fourth"
94
95 str(participant)
96 participant$level <- as.factor(participant$level)
97 levels(participant$level) = c("first", "second", "third", "fourth")
98
99 participant$group <- as.numeric(participant$level)
100
101 participant
102
103 # Exercise 1.
104 #####
105 # Print out bulit-in R data frame
106 mtcars
107
108 # Try the function head() on mtcars.
109 head(mtcars)           # head() enables us to see the first few observations of a data
110                         # frame.
111
112 # Similarly tail() prints out the last few observation in the data set.
113 tail(mtcars)
114
115 str(mtcars)
116
117 # Regress mpg on cyl and hp. (Chekc lm() function)
118
119 # Select the rows for mercedes. Hint: They include "Merc" and check grep() function
120
121 # Exercise 2.
122 #####
123 name <- c("Apple", "MS", "Google", "Honda", "GM", "Volks", "Hyundai", "Amazon")
124 type <- c("IT", "IT", "IT", "Auto", "Auto", "Auto", "Auto", "IT")
125 stock <- c(165.5, 55.48, 1119.20, 36.16, 41, 172.06, 162.5, 1429.95)
126 US <- c(TRUE, TRUE, TRUE, FALSE, TRUE, FALSE, FALSE, TRUE)
127
128 # Creat a data from from the vectors
129 portfolio <-
130 rm(name, type, stock, US)

```

```
131
132 # Check the structure of planet_df (Use str() function)
133
134 portfolio$name <- as.character(portfolio$name)
135
136 # Selection from data.frame.
137 # Print out stock price of google
138
139 # Print out data for Google (entire third row)
140
141 # Print out the first 5 values of stock column.
142
143 # Print out data for IT companies.
144
145 # Print out data for companies whose stock price is lower than Apple.
146
147
148 # Sort data.frame using order() function. From the most expensive one.
149 # Use rank to sort portforlio
150
151
```