```r
1    # Data structure: 3. Factor
2
3    # Factor a type of data structures used to store categorical variables. (e.g., gender)
4
5    # Gender vector
6    gender <- c("Male", "Female", "Female", "Male", "Male")
7    is.vector(gender)
8    class(gender)
9
10   # Convert gender_vector to a factor
11   factor_gender <- factor(gender)
12   factor_gender
13
14   # There are two different types of categorical variables:
15
16   # nominal categorical variable and ordinal categorical variable.
17   # 1. Norminal categorical variable: There is no implied order among categories.
18   # For example, Male and Female, and Cat, Dog, and Turtle
19   # 2. Ordinal categorical variable: There is a natural ordering.
20   # For example, "Low", "Medium" and "High", and "Primary", "Middle", "High"
21
22   # No ordering
23   pet <- c("Cat", "Dog", "Turtle", "Dog", "Cat", "Cat")
24   factor_pet <- factor(pet)
25   factor_pet
26   class(factor_pet)
27   factor_pet[1] > factor_pet[2]
28
29   # Natural ordering
30   income <- c("High", "High", "Low", "Midium", "Low")
31   factor_income <- factor(income, order = TRUE, levels = c("Low", "Midium", "High")) #
     You can specify the order (or level)
32   factor_income
33   class(factor_income)
34
35   income[3] > income[1]
36   factor_income[3] > factor_income[1]
37
38   # In survey, abbreviations are often used because it is convenient to record.
39   # But this can be confusing when you use survey data.
40   # You can recover the full words using levels function.
41   edu <- c("E", "E", "H", "M", "H", "C")
42   factor_edu <- factor(edu, order = TRUE, levels <- c("E", "M", "H", "C")) # If you don't
     specify the level, it will be alphabetical.
43   levels(factor_edu) <- c("Elemenary","Middle","High", "College")
44   factor_edu
45   as.numeric(factor_edu)
46
47   # summarize the factor
48   summary(factor_edu)
49   table(factor_edu)
50   summary(edu)
51
52   # Exercise 1.
     ################################################################################
53
54   set.seed(pi)
55   r <- rnorm(2754, 0, 1)
56   income <- exp(r)
57   hist(income, breaks=100)
58
59   # find 0.25, 0.5, 0.75, 0.95 quantiles of income. You can use quantile.
60   quant <- quantile(income, c(0.25, 0.5, 0.75, 0.95))
61
62   # Construct a vector, income.level, as follows:
63   # If income is <= 0.25 quantile, "VL"
64   # If income is > 0.25 quantile and <= 0.5 quantile, "L"
65   # If income is > 0.5 quantile and <= 0.75 quantile, "M"
66   # If income is > 0.75 quantile and <= 0.95 quantile, "H"
```

```r
67    # If income is > 0.95 quantile, "VH"
68
69
70    # Make an ordered factor from income.level. Specify the levels as
      c("VL","L","M","H","VH")
71    factor_income.level <-
72
73    levels(factor_income.level) <- c("Very Low", "Low", "Middle", "High", "Very High")
74    summary(factor_income.level)      # You can see the summary of factor_income.level.
75
76    # Construct a subvector income.high that includes income that belongs to "High" and
      "Very High".
77    income.high <-
78    hist(income.high, breaks = 20)
79
80    # Calculate the average income of people who belong to "Middle" and "High".
81
82
83    # What is the difference between average income of "Very High" and average income of
      "High"
84
85
86
87    # Exercise 2.
      #############################################################################
88
89    industry <- sample(c("Manufacture", "Service", "IT"), 100, replace=TRUE, prob=c(0.3,
      0.5, 0.2))
90    stock <- rep(NA,100)
91    stock[industry == "Manufacture"] <- rnorm(sum(industry=="Manufacture"), 3, 2)
92    stock[industry == "Service"] <- rnorm(sum(industry=="Service"), 2, 4)
93    stock[industry == "IT"] <- rnorm(sum(industry=="IT"), 8, 8)
94
95    factor.industry <- factor(industry)
96
97    # How many manufacturing, service, and IT companies?
98
99
100   # Compare the average stock prices and their standard deviations among these three
      industries.
101   mean(stock[factor.industry == "Manufacture"])
102   mean(stock[factor.industry == "Service"])
103   mean(stock[factor.industry == "IT"])
104
105
106
```