

RWORKSHEET #7a

Ann Erika D. Gabales

2022-12-23

1. Create a data frame for the table below.

```
df <- data.frame(Student = seq(1:10),  
                  PreTest = c(55,54,47,57,51,61,57,54,63,58),  
                  PostTest = c(61,60,56,63,56,63,59,56,62,61))  
df
```

##	Student	PreTest	PostTest
## 1	1	55	61
## 2	2	54	60
## 3	3	47	56
## 4	4	57	63
## 5	5	51	56
## 6	6	61	63
## 7	7	57	59
## 8	8	54	56
## 9	9	63	62
## 10	10	58	61

a. Compute the descriptive statistics using different packages (Hmisc and pastecs). Write the codes and its result.

```
library(Hmisc)
```

```
## Warning: package 'Hmisc' was built under R version 4.2.2
```

```
## Loading required package: lattice
```

```
## Loading required package: survival
```

```
## Loading required package: Formula
```

```
## Loading required package: ggplot2
```

```
## Warning: package 'ggplot2' was built under R version 4.2.2
```

```
##
## Attaching package: 'Hmisc'
```

```
## The following objects are masked from 'package:base':
##
##     format.pval, units
```

```
library(pastecs)
```

```
## Warning: package 'pastecs' was built under R version 4.2.2
```

```
describe(df)
```

```
## df
##
## 3 Variables      10 Observations
## -----
## Student
##      n missing distinct      Info      Mean      Gmd      .05      .10
##      10      0      10      1      5.5      3.667      1.45      1.90
##      .25      .50      .75      .90      .95
##      3.25      5.50      7.75      9.10      9.55
##
## lowest : 1 2 3 4 5, highest: 6 7 8 9 10
##
## Value      1 2 3 4 5 6 7 8 9 10
## Frequency  1 1 1 1 1 1 1 1 1 1
## Proportion 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1
## -----
## PreTest
##      n missing distinct      Info      Mean      Gmd
##      10      0      8      0.988      55.7      5.444
##
## lowest : 47 51 54 55 57, highest: 55 57 58 61 63
##
## Value      47 51 54 55 57 58 61 63
## Frequency  1 1 2 1 2 1 1 1
## Proportion 0.1 0.1 0.2 0.1 0.2 0.1 0.1 0.1
## -----
## PostTest
##      n missing distinct      Info      Mean      Gmd
##      10      0      6      0.964      59.7      3.311
##
## lowest : 56 59 60 61 62, highest: 59 60 61 62 63
##
## Value      56 59 60 61 62 63
## Frequency  3 1 1 2 1 2
## Proportion 0.3 0.1 0.1 0.2 0.1 0.2
## -----
```

```
stat.desc(df)
```

```
##           Student      PreTest      PostTest
## nbr.val      10.0000000  10.0000000  10.0000000
## nbr.null      0.0000000   0.0000000   0.0000000
## nbr.na        0.0000000   0.0000000   0.0000000
## min           1.0000000  47.0000000  56.0000000
## max           10.0000000  63.0000000  63.0000000
## range         9.0000000  16.0000000   7.0000000
## sum           55.0000000 557.0000000 597.0000000
## median        5.5000000  56.0000000  60.5000000
## mean          5.5000000  55.7000000  59.7000000
## SE.mean       0.9574271   1.46855938  0.89504811
## CI.mean.0.95  2.1658506   3.32211213  2.02473948
## var           9.1666667  21.56666667  8.01111111
## std.dev       3.0276504   4.64399254  2.83039063
## coef.var      0.5504819   0.08337509  0.04741023
```

2.The Department of Agriculture was studying the effects of several levels of a fertilizer on the growth of a plant. For some analyses, it might be useful to convert the fertilizer levels to an ordered factor.

```
DepartmentofAgriculture <- c(10,10,10,20,20,50,10,
                             20,10,50,20,50,20,10)
```

a. Write the codes and describe the result.

```
ordered_factor <- sort(DepartmentofAgriculture, decreasing = FALSE)
ordered_factor
```

```
## [1] 10 10 10 10 10 10 20 20 20 20 20 50 50 50
```

As what expected, the result was sorted in increasing number (from 10 to 50).

3. Abdul Hassan, president of Floor Coverings Unlimited, has asked you to study the exercise levels undertaken by 10 subjects were “l”, “n”, “n”, “i”, “l”, “l”, “n”, “n”, “i”, “l” ; n=none, l=light, i=intense

###a. What is the best way to represent this in R?

```
#DATAFRAME
Subjects <- c("l","n","n","i","l","l","n","n","i","l")

df <- data.frame(Subjects)
df
```

```
##      Subjects
## 1      1
## 2      n
## 3      n
## 4      i
## 5      1
## 6      1
## 7      n
## 8      n
## 9      i
## 10     1
```

4. Sample of 30 tax accountants from all the states and territories of Australia and their individual state of origin is specified by a character vector of state mnemonics as:

```
state <- c("tas", "sa", "qld", "nsw", "nsw", "nt", "wa", "wa", "qld",
           "vic", "nsw", "vic", "qld", "qld", "sa", "tas", "sa", "nt",
           "wa", "vic", "qld", "nsw", "nsw", "wa", "sa", "act", "nsw",
           "vic", "vic", "act")
```

a. Apply the factor function and factor level. Describe the results.

```
fcFunction <- factor(state)
fcFunction
```

```
## [1] tas sa qld nsw nsw nt wa wa qld vic nsw vic qld qld sa tas sa nt wa
## [20] vic qld nsw nsw wa sa act nsw vic vic act
## Levels: act nsw nt qld sa tas vic wa
```

the vectors were factored and it also has been leveled. ## 5. From #4 - continuation: -Suppose we have the incomes of the same tax accountants in another vector (insuitably large units of money)

```
incomes <- c(60, 49, 40, 61, 64, 60, 59, 54,
             62, 69, 70, 42, 56, 61, 61, 61, 58, 51, 48,
             65, 49, 49, 41, 48, 52, 46, 59, 46, 58, 43)
```

a. Calculate the sample mean income for each state we can now use the special function `tapply()`:

```
incmeans <- tapply(incomes, state, mean)
incmeans
```

```
##      act      nsw      nt      qld      sa      tas      vic      wa
## 44.50000 57.33333 55.50000 53.60000 55.00000 60.50000 56.00000 52.25000
```

b. Copy the results and interpret.

As what shown on the result, each has its own sample mean income.

6. Calculate the standard errors of the state income means (refer again to number 3)

a. What is the standard error? Write the codes.

```
stdError <- function(x) sqrt(var(x)/length(x))
incster <- tapply(incomes, state, stdError)
incster
```

```
##      act      nsw      nt      qld      sa      tas      vic      wa
## 1.500000 4.310195 4.500000 4.106093 2.738613 0.500000 5.244044 2.657536
```

b. Interpret the result.

As what shown on the result, each state income means has standard errors.

7. Use the titanic dataset.

```
data("Titanic")
dfTitanic<- data.frame(Titanic)
```

a. subset the titatic dataset of those who survived and not survived. Show the codes and its result.

```
headSubset <- subset(dfTitanic, select = "Survived")
headSubset
```

```
##      Survived
## 1          No
## 2          No
## 3          No
## 4          No
## 5          No
## 6          No
## 7          No
## 8          No
## 9          No
## 10         No
## 11         No
## 12         No
## 13         No
## 14         No
```

```
## 15      No
## 16      No
## 17     Yes
## 18     Yes
## 19     Yes
## 20     Yes
## 21     Yes
## 22     Yes
## 23     Yes
## 24     Yes
## 25     Yes
## 26     Yes
## 27     Yes
## 28     Yes
## 29     Yes
## 30     Yes
## 31     Yes
## 32     Yes
```

8. The data sets are about the breast cancer Wisconsin. The samples arrive periodically as Dr. Wolberg reports his clinical cases. The database therefore reflects this chronological grouping of the data. You can create this dataset in Microsoft Excel.

a. describe what is the dataset all about.

The dataset is all about the breast cancer Wisconsin whereas it can be seen the brief informations about the respondents who were diagnosed of Breast cancer.

b. Import the data from MS Excel. Copy the codes.

```
library(readxl)
```

```
## Warning: package 'readxl' was built under R version 4.2.2
```

```
Breast_Cancer <- read_excel("C:/Users/asus/Desktop/GABALES_BSIT2A/Breast_Cancer.xlsx")
View(Breast_Cancer)
```

c. Compute the descriptive statistics using different packages. Find the values of:

c.1 Standard error of the mean for clump thickness.

```
stdCLth <- stdError(Breast_Cancer$`CL. thickness`)
stdCLth
```

```
## [1] 0.4092884
```

c.2 Coefficient of variability for Marginal Adhesion.

```
COV <- sd(Breast_Cancer$`Marg. Adhesion`) / mean(Breast_Cancer$`Marg. Adhesion`)* 100  
COV
```

```
## [1] 97.67235
```

c.3 Number of null values of Bare Nuclei.

```
Null_Values <- subset(Breast_Cancer[ , 7], `Bare. Nuclei` == "NA")  
Null_Values
```

```
## # A tibble: 2 x 1  
##   'Bare. Nuclei'  
##   <chr>  
## 1 NA  
## 2 NA
```

c.4 Mean and standard deviation for Bland Chromatin

```
mean(Breast_Cancer$`Bl. Chromatin`)
```

```
## [1] 3.836735
```

```
sd(Breast_Cancer$`Bl. Chromatin`)
```

```
## [1] 2.085135
```

c.5 Confidence interval of the mean for Uniformity of Cell Shape

```
CalMean<- mean(Breast_Cancer$`Cell Shape`)  
CalMean
```

```
## [1] 3.163265
```

d. How many attributes?

```
length(attributes(Breast_Cancer))
```

```
## [1] 3
```

e. Find the percentage of respondents who are malignant. Interpret the results.

```
P_respond <- subset(Breast_Cancer, Class == "malignant")
P_respond
```

```
## # A tibble: 17 x 11
##       ID CL. thickne~1 Cell ~2 Cell ~3 Marg.~4 Epith~5 Bare.~6 Bl. C~7 Norma~8
##       <dbl>         <dbl>   <dbl>   <dbl>   <dbl>   <dbl> <chr>     <dbl>   <dbl>
##  1 1041801           5         3         3         3         2 3         4         4
##  2 1044572           8         7         5        10         7 9         5         5
##  3 1047630           7         4         6         4         6 1         4         3
##  4 1050670          10         7         7         6         4 10        4         1
##  5 1054590           7         3         2        10         5 10        5         4
##  6 1054593          10         5         5         3         6 7         7        10
##  7 1057013           8         4         5         1         2 NA         7         3
##  8 1065726           5         2         3         4         2 7         3         6
##  9 1072179          10         7         7         3         8 5         7         4
## 10 1080185          10        10        10         8         6 1         8         9
## 11 1084584           5         4         4         9         2 10        5         6
## 12 1091262           2         5         3         3         6 7         7         5
## 13 1099510          10         4         3         1         3 3         6         5
## 14 1100524           6        10        10         2         8 10        7         3
## 15 1102573           5         6         5         6        10 1         3         1
## 16 1103608          10        10        10         4         8 1         8        10
## 17 1105257           3         7         7         4         4 9         4         8
## # ... with 2 more variables: Mitoses <dbl>, Class <chr>, and abbreviated
## #   variable names 1: 'CL. thickness', 2: 'Cell size', 3: 'Cell Shape',
## #   4: 'Marg. Adhesion', 5: 'Epith. C.size', 6: 'Bare. Nuclei',
## #   7: 'Bl. Cromatin', 8: 'Normal nucleoli'
```

```
library(scales)
```

```
## Warning: package 'scales' was built under R version 4.2.2
```

```
percent_respo <- 17 / 49
percent(percent_respo)
```

```
## [1] "35%"
```

9. Export the data abalone to the Microsoft excel file. Copy the codes.

```
library("AppliedPredictiveModeling")
```

```
## Warning: package 'AppliedPredictiveModeling' was built under R version 4.2.2
```

```
data("abalone")
head(abalone)
```

```
##   Type LongestShell Diameter Height WholeWeight ShuckedWeight VisceraWeight
## 1    M      0.455     0.365  0.095      0.5140      0.2245      0.1010
```



```
## 2      M      0.350    0.265 0.090      0.2255      0.0995      0.0485
## 3      F      0.530    0.420 0.135      0.6770      0.2565      0.1415
## 4      M      0.440    0.365 0.125      0.5160      0.2155      0.1140
## 5      I      0.330    0.255 0.080      0.2050      0.0895      0.0395
## 6      I      0.425    0.300 0.095      0.3515      0.1410      0.0775
##      ShellWeight Rings
## 1      0.150     15
## 2      0.070      7
## 3      0.210      9
## 4      0.155     10
## 5      0.055      7
## 6      0.120      8
```

```
summary(abalone)
```

```
##      Type      LongestShell      Diameter      Height      WholeWeight
## F:1307  Min.    :0.075    Min.    :0.0550  Min.    :0.0000  Min.    :0.0020
## I:1342  1st Qu.:0.450    1st Qu.:0.3500  1st Qu.:0.1150  1st Qu.:0.4415
## M:1528  Median :0.545    Median :0.4250  Median :0.1400  Median :0.7995
##          Mean   :0.524    Mean   :0.4079  Mean   :0.1395  Mean   :0.8287
##          3rd Qu.:0.615    3rd Qu.:0.4800  3rd Qu.:0.1650  3rd Qu.:1.1530
##          Max.   :0.815    Max.   :0.6500  Max.   :1.1300  Max.   :2.8255
## ShuckedWeight  VisceraWeight  ShellWeight  Rings
## Min.    :0.0010  Min.    :0.0005  Min.    :0.0015  Min.    : 1.000
## 1st Qu.:0.1860  1st Qu.:0.0935  1st Qu.:0.1300  1st Qu.: 8.000
## Median :0.3360  Median :0.1710  Median :0.2340  Median : 9.000
## Mean   :0.3594  Mean   :0.1806  Mean   :0.2388  Mean   : 9.934
## 3rd Qu.:0.5020  3rd Qu.:0.2530  3rd Qu.:0.3290  3rd Qu.:11.000
## Max.   :1.4880  Max.   :0.7600  Max.   :1.0050  Max.   :29.000
```

```
dfabalone <- data.frame(abalone)
```

```
#Exporting the data abalone to the Microsoft excel file
```

```
library(xlsx)
```

```
## Warning: package 'xlsx' was built under R version 4.2.2
```

```
write.xlsx(dfabalone, "C:/Users/asus/Desktop/GABALES_BSIT2A/Abalone.xlsx", col.names = TRUE)
```