

# Worksheet #7a

Quennie Tabladillo

2022-12-12

##R Markdown

1. Create a data frame for the table below.

```
Student_Scores <- data.frame(  
  Student = c(1:10),  
  Pre_test = c(55,54,47,57,51,61,57,54,63,58),  
  Post_test = c(61,60,56,63,56,63,59,56,62,61)  
)  
Student_Scores
```

##	Student	Pre_test	Post_test
## 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(Student_Scores)
```

```
## Student_Scores
```

```
##
## 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

```
## -----
```

```
## Pre_test
```

	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

```
## -----
```

```
## Post_test
```

	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(Student_Scores)
```

```
##           Student      Pre_test      Post_test
## nbr.val      10.0000000  10.00000000  10.00000000
## nbr.null      0.0000000  0.00000000  0.00000000
## nbr.na        0.0000000  0.00000000  0.00000000
## min           1.0000000  47.00000000  56.00000000
## max          10.0000000  63.00000000  63.00000000
## range         9.0000000  16.00000000  7.00000000
## sum          55.0000000 557.00000000 597.00000000
## median        5.5000000  56.00000000  60.50000000
## mean          5.5000000  55.70000000  59.70000000
## 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.

a. the data were 10,10,10, 20,20,50,10,20,10,50,20,50,20,10.

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

fertilizer_factor <- factor(fertilizer, ordered = TRUE)
fertilizer_factor
```

```
## [1] 10 10 10 20 20 50 10 20 10 50 20 50 20 10
## Levels: 10 < 20 < 50
```

3. Abdul Hassan, president of Floor Coverings Unlimited, has asked you to study the exercise levels undertaken by 10 num3 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?

```
exercise_levels <- c("l","n","n","i","l","l","n","n","i","l")

data.frame(exercise_levels)
```

```
##      exercise_levels
## 1                   l
## 2                   n
## 3                   n
## 4                   i
## 5                   l
## 6                   l
```

```
## 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.

```
australia <- factor(state)
australia
```

```
## [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
```

5. From #4 - continuation: Suppose we have the incomes of the same tax accountants in another vector (in suitably 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()`:

```
sample_5 <- tapply(incomes, state, mean)
sample_5
```

```
##      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

```
#results
#act      nsw      nt      qld      sa      tas      vic      wa
#44.50000 57.33333 55.50000 53.60000 55.00000 60.50000 56.00000 52.25000
#the results shows the mean of every level of income for each state.
```

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

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

```

standard_err.n <- length(sample_5)
standard_err.sd <- sd(sample_5)
standard_err.se <- standard_err.sd/sqrt(standard_err.n)
standard_err.se

```

```
## [1] 1.653911
```

b. Interpret the result

*#This is how I get the state income means by dividing the sd() to sqrt() or length()  
#and that is how I get the standard errors of the state income means and this was  
#the result.*

7. Use the titanic dataset

a.Subset the titanic dataset of those who survived and not survived. Show the codes and its result.

```

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

survived <- subset(Titanic, Survived == "Yes")
survived

```

```

##      Class      Sex   Age Survived Freq
## 17   1st    Male Child      Yes     5
## 18   2nd    Male Child      Yes    11
## 19   3rd    Male Child      Yes    13
## 20  Crew    Male Child      Yes     0
## 21   1st Female Child      Yes     1
## 22   2nd Female Child      Yes    13
## 23   3rd Female Child      Yes    14
## 24  Crew Female Child      Yes     0
## 25   1st    Male Adult      Yes    57
## 26   2nd    Male Adult      Yes    14
## 27   3rd    Male Adult      Yes    75
## 28  Crew    Male Adult      Yes   192
## 29   1st Female Adult      Yes   140
## 30   2nd Female Adult      Yes    80
## 31   3rd Female Adult      Yes    76
## 32  Crew Female Adult      Yes    20

```

```

not_survived <- subset(Titanic, Survived == "No")
not_survived

```

```

##      Class      Sex   Age Survived Freq
## 1   1st    Male Child      No     0
## 2   2nd    Male Child      No     0
## 3   3rd    Male Child      No    35
## 4  Crew    Male Child      No     0
## 5   1st Female Child      No     0

```

```
## 6    2nd Female Child      No    0
## 7    3rd Female Child      No   17
## 8    Crew Female Child      No    0
## 9    1st   Male Adult      No  118
## 10   2nd   Male Adult      No  154
## 11   3rd   Male Adult      No  387
## 12   Crew   Male Adult      No  670
## 13   1st Female Adult      No    4
## 14   2nd Female Adult      No   13
## 15   3rd Female Adult      No   89
## 16   Crew Female Adult      No    3
```

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.

```
#Answer: The dataset is all about the breast cancer Wisconsin.
```

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("E:/CS 101/Worksheet 7a/Breast_Cancer.xlsx")
```

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

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

```
Breast_cancerc1.n <- length(Breast_cancer$'CL. thickness')
Breast_cancerc1.sd <- sd(Breast_cancer$'CL. thickness')
Breast_cancerc1.se <- Breast_cancerc1.sd/sqrt(Breast_cancer$'CL. thickness')
Breast_cancerc1.se
```

```
## [1] 1.2812754 1.2812754 1.6541194 1.1696391 1.4325095 1.0129371 2.8650189
## [8] 2.0258743 2.0258743 1.4325095 2.8650189 2.0258743 1.2812754 2.8650189
## [15] 1.0129371 1.0828754 1.4325095 1.4325095 0.9059985 1.1696391 1.0828754
## [22] 0.9059985 1.6541194 1.0129371 2.8650189 1.2812754 1.6541194 1.2812754
## [29] 2.0258743 2.8650189 1.6541194 2.0258743 0.9059985 2.0258743 1.6541194
## [36] 2.0258743 0.9059985 1.1696391 1.2812754 2.0258743 1.1696391 0.9059985
## [43] 1.1696391 1.2812754 0.9059985 2.8650189 1.6541194 2.8650189 1.4325095
```

c.2 Coefficient of variability for Marginal Adhesion.

```
sd(Breast_cancer$`Marg. Adhesion`) / mean(Breast_cancer$`Marg. Adhesion`) * 100
```

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

c.3 Number of null values of Bare Nuclei

```
Breast_cancerc3 <- subset(Breast_cancer, `Bare. Nuclei` == "NA")
Breast_cancerc3
```

```
## # A tibble: 2 x 11
##       Id CL. t~1 Cell ~2 Cell ~3 Marg.~4 Epith~5 Bare.~6 Bl. C~7 Norma~8 Mitoses
##       <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl> <chr>       <dbl>   <dbl>   <dbl>
## 1 1.06e6      8      4      5      1      2 NA          7      3      1
## 2 1.10e6      6      6      6      9      6 NA          7      8      1
## # ... with 1 more variable: 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'
```

c.4 Mean and standard deviation for Bland Chromatin

```
mean(Breast_cancer$`Bl. Cromatin`)
```

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

```
sd(Breast_cancer$`Bl. Cromatin`)
```

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

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

```
breast_cancerc5 <- mean(Breast_cancer$`Cell Shape`)
breast_cancerc5
```

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

```
#Calculate the standard error of the mean
numA <- length(Breast_cancer$`Cell Shape`)
numB <- sd(Breast_cancer$`Cell Shape`)
numC <- numB/sqrt(numA)
numC
```

```
## [1] 0.4158294
```

```
#Find the t-score that corresponds to the confidence level
numD = 0.05
numE = numA - 1
numF = qt(p=numD/2, df=numE, lower.tail=F)
numF
```

```
## [1] 2.010635
```

```
#Constructing the confidence interval
numG <- numF * numC
numG
```

```
## [1] 0.836081
```

```
#Lower
numH <- breast_cancerc5 - numG
numH
```

```
## [1] 2.327184
```

```
#Upper
numI <- breast_cancerc5 + numG
numI
```

```
## [1] 3.999346
```

```
c(numH,numI)
```

```
## [1] 2.327184 3.999346
```

d. How many attributes?

```
attributes(Breast_cancer)

## $class
## [1] "tbl_df"      "tbl"        "data.frame"
##
## $row.names
## [1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
## [26] 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49
##
## $names
## [1] "Id"                "CL. thickness"    "Cell size"        "Cell Shape"
## [5] "Marg. Adhesion"    "Epith. C.size"    "Bare. Nuclei"      "Bl. Cromatin"
## [9] "Normal nucleoli"   "Mitoses"          "Class"
```

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

```
cancer <- subset(Breast_cancer, Class == "malignant")
cancer

## # A tibble: 1 x 11
##       Id CL. t~1 Cell ~2 Cell ~3 Marg.~4 Epith~5 Bare.~6 Bl. C~7 Norma~8 Mitoses
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <chr> <dbl> <dbl> <dbl>
## 1 1.02e6      8     10     10      8      7 10      9      7      1
## # ... with 1 more variable: 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'
```

```
17 / 49 * 100
```

```
## [1] 34.69388
```



```
#There are 34.69388 or 35% of respondents who are malignant.
```

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
```

```
#exporting the data abalone
```

```
library(xlsx)
```

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

```
write.xlsx("abalone", "C:\\Users\\Quennie\\Documents\\abalone.xlsx")
```