

# RWorksheet\_5

Jiruel Suero BSIT 2-C

2023-12-14

1. Create a data frame for the table below. Show your solution.
  - a. Compute the descriptive statistics using different packages (Hmisc and pastecs). Write the codes and its result.

```
StudentScore <- data.frame(Student = c(1, 2, 3, 4, 5, 6, 7, 8, 9, 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))
```

StudentScore

##	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

```
library(Hmisc)
```

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

```
library(pastecs)
```

```
HmiscStats <- describe(StudentScore[, c("PreTest", "PostTest")])
HmiscStats
```

```
## StudentScore[, c("PreTest", "PostTest")]
```

```
##
```

```
## 2 Variables      10 Observations
```

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

```
## PreTest
```

##	n	missing	distinct	Info	Mean	Gmd
##	10	0	8	0.988	55.7	5.444

```
##
```

```
## 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
##
## For the frequency table, variable is rounded to the nearest 0
## -----
## PostTest
##      n missing distinct      Info      Mean      Gmd
##      10      0        6     0.964     59.7     3.311
##
## 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
##
## For the frequency table, variable is rounded to the nearest 0
## -----

pastecsStats <- stat.desc(StudentScore[, c('PreTest', 'PostTest')])
pastecsStats
```

```
##              PreTest      PostTest
## nbr.val      10.00000000 10.00000000
## nbr.null      0.00000000  0.00000000
## nbr.na        0.00000000  0.00000000
## min          47.00000000 56.00000000
## max          63.00000000 63.00000000
## range        16.00000000  7.00000000
## sum          557.00000000 597.00000000
## median        56.00000000 60.50000000
## mean          55.70000000 59.70000000
## SE.mean       1.46855938  0.89504811
## CI.mean.0.95  3.32211213  2.02473948
## var          21.56666667  8.01111111
## std.dev       4.64399254  2.83039063
## coef.var      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.

```
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
##
## The following objects are masked from 'package:stats':
##
##     first, last
##
## The following objects are masked from 'package:Hmisc':
##
##     src, summarize
##
## The following objects are masked from 'package:base':
##
##     filter, lag
##
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

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

orderedFactor <- factor(fertilizerLevels, levels = unique(fertilizerLevels))

basicStats <- summary(orderedFactor)
basicStats

## 10 20 50
##  6  5  3
```

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?

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

ExerciseFactor <- factor(exercerciseLevels, levels = c("n","l","i"))

basic_stats <- summary(ExerciseFactor)
basic_stats

## n l i
## 5 4 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.

```
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")
stateFactor <- factor(state)
stateFactor

## [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

summaryState <- summary(stateFactor)
```

*#The output will show the levels (unique values) in the factor (act, nsw, nt, qld, sa, tas, vic, wa) and*

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():

```
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)
```

```
meanIncome <- tapply(incomes, stateFactor, mean)
meanIncome
```

```
##      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.

```
# 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 code attempts to calculate the mean income for different states using the tapply function, but it*

6. Calculate the standard errors of the state income means (refer again to number 3) `stdError <- function(x) sqrt(var(x)/length(x))` Note: After this assignment, the standard errors are calculated by: `incster <- tapply(incomes, statef, stdError)` a. What is the standard error? Write the codes.

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

```
##      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.

*#These values indicate the precision of the estimated mean for each region. Higher standard errors gene*

7. Use the titanic dataset.

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

```
install.packages("titanic")
```

```
## Installing package into '/cloud/lib/x86_64-pc-linux-gnu-library/4.3'
## (as 'lib' is unspecified)
```

```
library(titanic)
```

```
data("titanic_train")
titanic_data <- titanic_train
```

```
survived_data <- subset(titanic_data, Survived == 1)
```

```
not_survived_data <- subset(titanic_data, Survived == 0)
```

```
head(survived_data)
```

```
##      PassengerId Survived Pclass
## 2              2         1       1
## 3              3         1       3
## 4              4         1       1
## 9              9         1       3
## 10             10         1       2
## 11             11         1       3
##
##              Name      Sex Age SibSp Parch
## 2 Cumings, Mrs. John Bradley (Florence Briggs Thayer) female 38      1      0
## 3                               Heikkinen, Miss. Laina female 26      0      0
## 4 Futrelle, Mrs. Jacques Heath (Lily May Peel) female 35      1      0
```

```
## 9      Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg) female 27      0      2
## 10      Nasser, Mrs. Nicholas (Adele Achem) female 14      1      0
## 11      Sandstrom, Miss. Marguerite Rut female 4      1      1
##      Ticket      Fare Cabin Embarked
## 2      PC 17599 71.2833      C85      C
## 3 STON/O2. 3101282 7.9250      S
## 4      113803 53.1000      C123      S
## 9      347742 11.1333      S
## 10     237736 30.0708      C
## 11     PP 9549 16.7000      G6      S
```

```
head(not_survived_data)
```

```
##      PassengerId Survived Pclass      Name Sex Age SibSp
## 1      1      0      3      Braund, Mr. Owen Harris male 22      1
## 5      5      0      3      Allen, Mr. William Henry male 35      0
## 6      6      0      3      Moran, Mr. James male NA      0
## 7      7      0      1      McCarthy, Mr. Timothy J male 54      0
## 8      8      0      3 Palsson, Master. Gosta Leonard male 2      3
## 13     13      0      3 Saundercock, Mr. William Henry male 20      0
##      Parch      Ticket      Fare Cabin Embarked
## 1      0 A/5 21171 7.2500      S
## 5      0 373450 8.0500      S
## 6      0 330877 8.4583      Q
## 7      0 17463 51.8625      E46      S
## 8      1 349909 21.0750      S
## 13     0 A/5. 2151 8.0500      S
```

```
survived_data <- titanic_data[titanic_data$Survived == 1, ]
```

```
not_survived_data <- titanic_data[titanic_data$Survived == 0, ]
```

```
head(survived_data)
```

```
##      PassengerId Survived Pclass
## 2      2      1      1
## 3      3      1      3
## 4      4      1      1
## 9      9      1      3
## 10     10     1      2
## 11     11     1      3
##      Name Sex Age SibSp Parch
## 2 Cumings, Mrs. John Bradley (Florence Briggs Thayer) female 38      1      0
## 3 Heikkinen, Miss. Laina female 26      0      0
## 4 Futrelle, Mrs. Jacques Heath (Lily May Peel) female 35      1      0
## 9 Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg) female 27      0      2
## 10 Nasser, Mrs. Nicholas (Adele Achem) female 14      1      0
## 11 Sandstrom, Miss. Marguerite Rut female 4      1      1
##      Ticket      Fare Cabin Embarked
## 2      PC 17599 71.2833      C85      C
## 3 STON/O2. 3101282 7.9250      S
## 4      113803 53.1000      C123      S
## 9      347742 11.1333      S
## 10     237736 30.0708      C
```

```
## 11          PP 9549 16.7000    G6          S
```

```
head(not_survived_data)
```

```
##      PassengerId Survived Pclass                    Name Sex Age SibSp
## 1              1         0      3      Braund, Mr. Owen Harris male  22     1
## 5              5         0      3      Allen, Mr. William Henry male  35     0
## 6              6         0      3              Moran, Mr. James male   NA     0
## 7              7         0      1      McCarthy, Mr. Timothy J male  54     0
## 8              8         0      3 Palsson, Master. Gosta Leonard male   2     3
## 13             13         0      3 Saunderson, Mr. William Henry male  20     0
##      Parch      Ticket    Fare Cabin Embarked
## 1         0 A/5 21171    7.2500         S
## 5         0  373450    8.0500         S
## 6         0  330877    8.4583         Q
## 7         0   17463   51.8625     E46     S
## 8         1  349909   21.0750         S
## 13        0 A/5. 2151    8.0500         S
```

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

chronologittps://drive.google.com/file/d/16MFLoehCgx2MJuNSAuB2CsBy6eDIr- u/view?usp=drive\_link)

a. describe what is the dataset all about.

*#The dataset consists of cytological features of breast cancer cell samples, such as clump thickness, s*

- d. Compute the descriptive statistics using different packages. Find the values of: d.1 Standard error of the mean for clump thickness.

```
library(readr)
```

```
breastcancer_wisconsin <- read_csv("/cloud/project/breastcancer_wisconsin.csv")
```

```
## Rows: 699 Columns: 11
## -- Column specification -----
## Delimiter: ","
## chr (1): bare_nucleoli
## dbl (10): id, clump_thickness, size_uniformity, shape_uniformity, marginal_a...
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
str(breastcancer_wisconsin)
```

```
## spc_tbl_ [699 x 11] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ id : num [1:699] 1000025 1002945 1015425 1016277 1017023 ...
## $ clump_thickness : num [1:699] 5 5 3 6 4 8 1 2 2 4 ...
## $ size_uniformity : num [1:699] 1 4 1 8 1 10 1 1 1 2 ...
## $ shape_uniformity : num [1:699] 1 4 1 8 1 10 1 2 1 1 ...
## $ marginal_adhesion: num [1:699] 1 5 1 1 3 8 1 1 1 1 ...
## $ epithelial_size : num [1:699] 2 7 2 3 2 7 2 2 2 2 ...
## $ bare_nucleoli : chr [1:699] "1" "10" "2" "4" ...
## $ bland_chromatin : num [1:699] 3 3 3 3 3 9 3 3 1 2 ...
## $ normal_nucleoli : num [1:699] 1 2 1 7 1 7 1 1 1 1 ...
## $ mitoses : num [1:699] 1 1 1 1 1 1 1 1 5 1 ...
## $ class : num [1:699] 2 2 2 2 2 4 2 2 2 2 ...
```

```
## - attr(*, "spec")=
## .. cols(
## ..   id = col_double(),
## ..   clump_thickness = col_double(),
## ..   size_uniformity = col_double(),
## ..   shape_uniformity = col_double(),
## ..   marginal_adhesion = col_double(),
## ..   epithelial_size = col_double(),
## ..   bare_nucleoli = col_character(),
## ..   bland_chromatin = col_double(),
## ..   normal_nucleoli = col_double(),
## ..   mitoses = col_double(),
## ..   class = col_double()
## .. )
## - attr(*, "problems")=<externalptr>
```

```
summary(breastcancer_wisconsin)
```

```
##           id           clump_thickness  size_uniformity  shape_uniformity
## Min.      : 61634   Min.      : 1.000   Min.      : 1.000   Min.      : 1.000
## 1st Qu.: 870688   1st Qu.: 2.000   1st Qu.: 1.000   1st Qu.: 1.000
## Median : 1171710   Median : 4.000   Median : 1.000   Median : 1.000
## Mean    : 1071704   Mean    : 4.418   Mean    : 3.134   Mean    : 3.207
## 3rd Qu.: 1238298   3rd Qu.: 6.000   3rd Qu.: 5.000   3rd Qu.: 5.000
## Max.    :13454352   Max.    :10.000   Max.    :10.000   Max.    :10.000
## marginal_adhesion epithelial_size  bare_nucleoli      bland_chromatin
## Min.      : 1.000   Min.      : 1.000   Length:699        Min.      : 1.000
## 1st Qu.: 1.000   1st Qu.: 2.000   Class :character   1st Qu.: 2.000
## Median : 1.000   Median : 2.000   Mode  :character   Median : 3.000
## Mean    : 2.807   Mean    : 3.216                      Mean    : 3.438
## 3rd Qu.: 4.000   3rd Qu.: 4.000                      3rd Qu.: 5.000
## Max.    :10.000   Max.    :10.000                      Max.    :10.000
## normal_nucleoli      mitoses              class
## Min.      : 1.000   Min.      : 1.000   Min.      :2.00
## 1st Qu.: 1.000   1st Qu.: 1.000   1st Qu.:2.00
## Median : 1.000   Median : 1.000   Median :2.00
## Mean    : 2.867   Mean    : 1.589   Mean    :2.69
## 3rd Qu.: 4.000   3rd Qu.: 1.000   3rd Qu.:4.00
## Max.    :10.000   Max.    :10.000   Max.    :4.00
```

d.2 Coefficient of variability for Marginal Adhesion.

```
colnames(breastcancer_wisconsin)
```

```
## [1] "id"           "clump_thickness" "size_uniformity"
## [4] "shape_uniformity" "marginal_adhesion" "epithelial_size"
## [7] "bare_nucleoli" "bland_chromatin" "normal_nucleoli"
## [10] "mitoses"      "class"
```

```
marginal_adhesion_cv <- sd(breastcancer_wisconsin$`Marginal Adhesion`) / mean(breastcancer_wisconsin$`M
```

```
## Warning: Unknown or uninitialised column: `Marginal Adhesion`.
```

```
## Unknown or uninitialised column: `Marginal Adhesion`.
```

```
## Warning in mean.default(breastcancer_wisconsin$`Marginal Adhesion`, na.rm =
```

```
## TRUE): argument is not numeric or logical: returning NA
```

```

marginal_adhesion_cv

## [1] NA

d.3 Number of null values of Bare Nuclei.
bare_nuclei_null_count <- sum(is.na(breastcancer_wisconsin$`Bare Nuclei`))

## Warning: Unknown or uninitialised column: `Bare Nuclei`.
bare_nuclei_null_count

## [1] 0

d.4 Mean and standard deviation for Bland Chromatin
clump_thickness_mean <- mean(breastcancer_wisconsin$clump_thickness)
clump_thickness_sd <- sd(breastcancer_wisconsin$clump_thickness)
clump_thickness_sem <- clump_thickness_sd / sqrt(length(breastcancer_wisconsin$clump_thickness))

clump_thickness_mean

## [1] 4.41774
clump_thickness_sd

## [1] 2.815741
clump_thickness_sem

## [1] 0.1065011

d.5 Confidence interval of the mean for Uniformity of Cell Shape
library(readr)

# Read the CSV file
data <- read_csv("/cloud/project/breastcancer_wisconsin.csv")

## Rows: 699 Columns: 11
## -- Column specification -----
## Delimiter: ","
## chr (1): bare_nucleoli
## dbl (10): id, clump_thickness, size_uniformity, shape_uniformity, marginal_a...
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
# Extract the column of interest
column_of_interest <- data$`Uniformity of Cell Shape`

## Warning: Unknown or uninitialised column: `Uniformity of Cell Shape`.
# Remove rows with missing values
column_of_interest_clean <- na.omit(column_of_interest)

# Calculate sample mean, sample size, and sample standard deviation using the cleaned data
sample_mean <- mean(column_of_interest_clean)

## Warning in mean.default(column_of_interest_clean): argument is not numeric or
## logical: returning NA

```



```

sample_size <- length(column_of_interest_clean)
sample_sd <- sd(column_of_interest_clean)

# Set the confidence level
confidence_level <- 0.95

# Calculate the margin of error using the t-distribution
margin_of_error <- qt((1 + confidence_level) / 2, df = sample_size - 1) * (sample_sd / sqrt(sample_size))

## Warning in qt((1 + confidence_level)/2, df = sample_size - 1): NaNs produced

# Calculate the confidence interval
confidence_interval <- c(sample_mean - margin_of_error, sample_mean + margin_of_error)

# Print the results
cat("Sample Mean:", sample_mean, "\n")

## Sample Mean: NA

cat("Confidence Interval:", confidence_interval[1], "to", confidence_interval[2], "\n")

```

## Confidence Interval: NA to NA

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

```

library("AppliedPredictiveModeling")
library(MASS)

```

```

##
## Attaching package: 'MASS'

## The following object is masked from 'package:dplyr':
##
##      select

```

```

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

```

```

str(abalone)

```

```

## 'data.frame':   4177 obs. of  9 variables:
##  $ Type          : Factor w/ 3 levels "F","I","M": 3 3 1 3 2 2 1 1 3 1 ...
##  $ LongestShell  : num  0.455 0.35 0.53 0.44 0.33 0.425 0.53 0.545 0.475 0.55 ...
##  $ Diameter      : num  0.365 0.265 0.42 0.365 0.255 0.3 0.415 0.425 0.37 0.44 ...

```

```
## $ Height      : num  0.095 0.09 0.135 0.125 0.08 0.095 0.15 0.125 0.125 0.15 ...
## $ WholeWeight : num  0.514 0.226 0.677 0.516 0.205 ...
## $ ShuckedWeight: num  0.2245 0.0995 0.2565 0.2155 0.0895 ...
## $ VisceraWeight: num  0.101 0.0485 0.1415 0.114 0.0395 ...
## $ ShellWeight  : num  0.15 0.07 0.21 0.155 0.055 0.12 0.33 0.26 0.165 0.32 ...
## $ Rings        : int  15 7 9 10 7 8 20 16 9 19 ...
```