We will introduce R in this lab. You should have installed R and Rstudio already.

Problem 1a: Manually Enter Data Into R.

The following data set represents the ages (in years) of a small class of students:

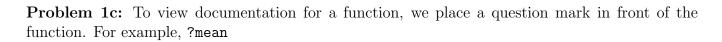
Enter the following command into the console:

where \leftarrow is an assignment operator and c() is the notation used in R to represent a vector of datapoints.

Problem 1b: R has functions that can quickly calculate summary statistics for us. Explore the following functions:

- mean()
- median()
- min()
- max()
- range()
- quantile()
- sd()
- var()

- summary()
- length()
- sum()
- IQR()
- table()



Choose one of the functions above that's not sd() and take a screenshot of its documentation. Explain how it works.

The max() function gives the maximum number in the age list.

Problem 1d: View the documentation for the sd() function. Does it divide by n-1 or n? As a result, does it calculate population standard deviation, or sample standard deviation?

It divides by n-1. As a result, it calculates sample standard deviation.

What if you have more than one variable recorded for each individual? For example, suppose the results of the first exam are recorded as follows:

$$\{75, 77, 94, 75, 79, 80, 66, 82, 86, 80, 78\}$$

Problem 2a: Enter this information into R by typing the following command into the R console (which stores the data into a variable named 'score'):

Problem 2b: To identify the data with each student, a student ID is assigned to each student. Type the following into the R console (which stores the identification numbers into a variable named 'id'):

$$id \leftarrow c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11)$$

You can also do: id <- 1:11

Problem 2c: It may be more convenient to create one object in R that contains all the information in one place. We use what is called a 'data frame' to do this. Type the following command into the R console:

Note: id = id means that we are going to name the column in the data frame 'id' (the name appearing to the left of the equal sign) and assign it the values stored in the vector named 'id' (the name appearing to the right of the equal sign). (CONT.)

Problem 2d: With the dataframe we just created, linformation from it. Explain what each command does	
• df[,1]	• df\$id
It gives out the first column of the data frame	It gives out the id column as a list from the data frame
• df[,2]	• df["id"]
It gives out the second column of the data frame	It extract the column called "id" from the data frame
• df[1,1]	• df\$id[3]
It gives out the first element of the first column in the data frame	It gives out the 3rd element of the column called "id" in the datafra
• df[2,2]	• df\$age
It gives out the second element of the second column in the data frame	It gives out the age column as a list from the data frame
• df[1,]	• df\$score
It gives out the first row of the data frame	It gives out the score column as a list from the data frame
• df[2,]	• df["score"]
It gives out the second row of the data frame	It extract the column called "id" from the data frame

Take a screenshot of your dataframe.

Now we learn how to read data from a Text File into R:

Problem 3a: First you need to set the working directory to where the file is located.

```
setwd("ENTER PATH HERE")
```

Problem 3b:

To read the data file into R and construct a data frame as done above, type the following command into the R console:

Note: If the file is a .csv file, use df <- read.csv("testData.csv")

Creating New Variables:

Problem 4a: Suppose we want to create a new variable age2 where we square each age value. We can do this by typing the following code into the R console:

If we want to add it to the data frame, we can type:

To view the data frame, enter the following command into the console: df

What do the following commands do?

It adds 1 to all the values in the age list

It adds each elements in the age list to itself(basically doubled the values)

Changing Data Values:

Problem 5a: Suppose the first age value was incorrectly recorded as 18. If the student's real age is 19, we can adjust the variable age using the following code:

If we want to make the change in the data frame df, the following will do the trick (pick one):

- df\$age[1] <- 19
- df[1,2] <- 19

Now, let's report summary statistics of age.

With the originally supplied age data:

$$\{18, 19, 23, 19, 24, 20, 18, 21, 22, 23, 18\}$$

fill out the table:

Statistic	Age
Count (n)	11
Mean	20.45455
St. Dev.	2.252272
Median	20
25%ile	18.5
75%ile	22.5
IQR	4
min	18
max	24
range	18~24

Describe the roles biostatistics serves in the discipline of public health.

To successfully complete this competency assessment, you may wish to review some of the following

material:

• From Whitlock: Chapter 1

• From Class: Lecture 1

Your written description should be double-space typed and under 250 words.

Biostatistics is a useful discipline helping scientists and health specialists testing hypothesis and

making decisions. It uses both quantitative and qualitative data to figure out practical applications

such as: if a drug is working on certain group of patients. The professor mentioned a famous

clinical trial about VitaminC in class. He mentioned that biostatistics analyzed if Vc works to

prevent

common cold. The scientist Linus Pauline tested the effect of Vc on volunteers along with a

control group. The results showed that his hypothesis was not statistically significant, which means

that Vc does not prevents common cold. He also mentioned his involvement in a project collecting

human milk to support pre-mature babies' survival, which implies the role of biostatistics includes

allocating resources and apply to patients needing the resources. It ensures new treatment/drugs

are safe to use and effective on patients.