

Lab 2: Data Structures

Complete all of the following questions, adding your inputs as code chunks (enclose within triple accent marks) within Rmarkdown.

The exercises are not marked and will not be factored into your course grade, but it is important to complete them to make sure you have the skills to answer assessment questions. You may consult any resource, including other students and the instructor. Please Knit this document to a PDF and upload your work via Canvas at the end of the session. Solutions will be posted for you to check your own answers.

String manipulation

1. Write your own sentence (at least 5 words) and store it in a string (character) variable.

```
sentence = 'This is an example sentence.'
```

2. Print out the first 10 characters of the string to the screen (use `substr`).

```
substr(sentence,1,10)
```

```
## [1] "This is an"
```

3. Print out the second character of the string.

```
substr(sentence,2,2)
```

```
## [1] "h"
```

4. Print the sentence in ALL UPPERCASE (use `toupper`).

```
toupper(sentence)
```

```
## [1] "THIS IS AN EXAMPLE SENTENCE."
```

5. Print out a string containing the sentence but with all spaces converted to underscores (use `gsub`).

```
gsub(' ', '_', sentence)
```

```
## [1] "This_is_an_example_sentence."
```

6. Create a vector where each element of the vector contains a word from the string (use `strsplit`) and print it out.

```
strsplit(sentence, ' ')[[1]]
```

```
## [1] "This"      "is"        "an"        "example"   "sentence."
```

7. Print the words in your string in alphabetical order by storing the vector above in a variable, and then sorting it with `sort`.

```
words = strsplit(sentence, ' ')[[1]]  
sort(words)
```

```
## [1] "an"        "example"   "is"        "sentence." "This"
```

8. Create a variable containing the string “34 miles”. (This is the approximate distance to Manchester.)

```
distancestr = '34 miles'
```

9. Use `strsplit` and `as.numeric` to extract the numerical portion of the string, convert it to a number, and then divide it by two in order to calculate half the distance to Manchester in miles without any new data input. Save it as a variable and print it to the screen.

```
halfdistance = as.numeric(strsplit(distancestr, ' ')[[1]][1])/2  
halfdistance
```

```
## [1] 17
```

10. Use the variable above to print out a single string that states: “Half the distance to Manchester is XX miles”, except with the true value in place of XX. (Use `paste` to attach the strings.)

```
paste('Half the distance to Manchester is',halfdistance,'miles')
```

```
## [1] "Half the distance to Manchester is 17 miles"
```

Matrices

11. Create a matrix (stored in variable `m`) of 3 rows and 5 columns, with every value equal to zero. Print out the matrix to confirm its shape.

```
m = matrix(0,3,5)  
m
```

```
##      [,1] [,2] [,3] [,4] [,5]  
## [1,]    0    0    0    0    0  
## [2,]    0    0    0    0    0  
## [3,]    0    0    0    0    0
```

12. Set the value in the 2nd row, 3rd column to 1, then print out the matrix again.

```
m[2,3] = 1  
m
```

```
##      [,1] [,2] [,3] [,4] [,5]  
## [1,]    0    0    0    0    0  
## [2,]    0    0    1    0    0  
## [3,]    0    0    0    0    0
```

13. Set all the values in the 4th column to 2. Set all values in the 3rd row to 3. Print out the matrix again.

```
m[,4] = 2  
m[3,] = 3  
m
```

```
##      [,1] [,2] [,3] [,4] [,5]  
## [1,]    0    0    0    2    0  
## [2,]    0    0    1    2    0  
## [3,]    3    3    3    3    3
```

14. Tack on a new row to the bottom of the matrix. The new row’s elements should be all equal to 4. Then, tack on a new column to the right of the matrix. The new column’s elements should be all equal to 5. (Use `rbind` and `cbind`.) Print out the matrix again and confirm that it now has 6 columns and 4 rows.

```
m = rbind(m,4)  
m = cbind(m,5)  
m
```

```
##      [,1] [,2] [,3] [,4] [,5] [,6]
## [1,]    0    0    0    2    0    5
## [2,]    0    0    1    2    0    5
## [3,]    3    3    3    3    3    5
## [4,]    4    4    4    4    4    5
```

15. Print the sum of all the elements in the matrix. (The result should be 60: if not, return to question 10).

```
sum(m)
```

```
## [1] 60
```

16. Print a vector containing the sums of all the elements in each individual row (use `apply`)

```
apply(m,1,sum)
```

```
## [1]  7  8 20 25
```

17. Print a vector containing the sums of all the elements in each individual column.

```
apply(m,2,sum)
```

```
## [1]  7  7  8 11  7 20
```

18. Print a vector containing the means of all the elements in each individual column.

```
apply(m,2,mean)
```

```
## [1] 1.75 1.75 2.00 2.75 1.75 5.00
```

19. Print a vector containing the medians of all the elements in each individual column.

```
apply(m,2,median)
```

```
## [1] 1.5 1.5 2.0 2.5 1.5 5.0
```

20. Print the transpose of the matrix (rows and columns swapped.)

```
t(m)
```

```
##      [,1] [,2] [,3] [,4]
## [1,]    0    0    3    4
## [2,]    0    0    3    4
## [3,]    0    1    3    4
## [4,]    2    2    3    4
## [5,]    0    0    3    4
## [6,]    5    5    5    5
```

Named vectors

21. Suppose you collected the following data on your friends' birthdays:

Alice: 31/8
 Bob: 5/5
 Carol: 11/12
 David: 17/10

Create a string vector containing *just* the birthdays.

```
bday = c('31/8', '5/5', '11/12', '17/10')
```

22. Add names to all four elements in your vector using your friends' names above.

```
names(bday) = c('Alice', 'Bob', 'Carol', 'David')
```

23. Print out Carol's birthday using her name as a subscript.

```
bday["Carol"]
```

```
##   Carol  
## "11/12"
```

Lists

Suppose you have some miscellaneous data about a car: It has 4 seats. Its top speed is 100 miles per hour. The make is Hyundai. The model is Elantra. The transmission is standard. The tyre pressure readings in PSI are 30, 31, 29, 33.

24. Create a list containing all this data in a single variable. Print out information about its structure.

```
carinfo = list(seats=4, topspeed=100, make='Hyundai', model='Elantra',  
              transmission='standard', tyrepressure=c(30,31,29,33))  
str(carinfo)
```

```
## List of 6  
## $ seats      : num 4  
## $ topspeed   : num 100  
## $ make       : chr "Hyundai"  
## $ model      : chr "Elantra"  
## $ transmission: chr "standard"  
## $ tyrepressure: num [1:4] 30 31 29 33
```

25. Print out the mean tyre pressure (using an operation on your list variable).

```
mean(carinfo$tyrepressure)
```

```
## [1] 30.75
```

Data Frames

26. Reload the student summary table from the in-class exercise.

```
survey = read.csv('survey.csv', as.is=FALSE)
```

27. Generate a summary of the dataframe using the **summary** command. Then, answer the following questions:

- What is the mean and median height of the students?
- What is the most popular football club?
- What is the most popular type of caffeinated beverage?

You can simply read them off the summary table and enter them down below the code chunk by hand. (In principle you can also calculate them using various R functions, but this is less straightforward.)

```
summary(survey)
```

```
##      age      height      gender      football_club  
## Min.   :21.00   Min.   :59.00  female:13   Liverpool   :8  
## 1st Qu.:22.00   1st Qu.:65.25   male  :19   Chelsea    :4
```

```
## Median :24.00 Median :68.00 none :2
## Mean :25.38 Mean :67.70 none :2
## 3rd Qu.:26.00 3rd Qu.:70.00 Real Madrid:2
## Max. :35.00 Max. :75.00 (Other) :6
## NA's :3 NA's :2 NA's :8
## beverage siblings av_height_est dogs hometown_pop
## coffee :19 Min. :0.000 Min. :59.00 No : 7 Min. : 300
## cola : 3 1st Qu.:1.000 1st Qu.:65.00 Yes:25 1st Qu.: 20500
## energy drink: 2 Median :1.000 Median :67.00 Median : 95000
## none : 1 Mean :1.774 Mean :66.40 Mean : 17309730
## tea : 5 3rd Qu.:2.500 3rd Qu.:69.75 3rd Qu.: 1540050
## NA's : 2 Max. :5.000 Max. :71.00 Max. :500000000
## NA's :1 NA's :2 NA's :1
## berlin_dist_est berlin_dist_unc distance fave_nums wake_time_wkday
## Min. : 100 Min. : 40.0 Min. : 0.000 3,7 : 2 07:00 :7
## 1st Qu.: 600 1st Qu.: 100.0 1st Qu.: 0.950 6,8 : 2 08:00 :7
## Median : 1000 Median : 200.0 Median : 1.750 10,46 : 1 07:30 :5
## Mean : 1616 Mean : 451.3 Mean : 2.928 13,21 : 1 06:00 :3
## 3rd Qu.: 2000 3rd Qu.: 500.0 3rd Qu.: 4.000 17,257 : 1 06:30 :3
## Max. :10000 Max. :5000.0 Max. :15.000 2,19 : 1 09:00 :3
## NA's :3 NA's :7 (Other):24 (Other):4
## wake_time_wkend colour1 colour2
## 07:00 : 5 Blue :14 Black:23
## 08:00 : 4 Green: 6 White: 9
## 08:30 : 4 Red :11
## 05:00 : 3 NA's : 1
## 09:00 : 3
## 10:00 : 3
## (Other):10
```

```
mean(survey$height,na.rm=TRUE); median(survey$height,na.rm=TRUE)
```

```
## [1] 67.7
```

```
## [1] 68
```

```
sort(table(survey$football_club),decreasing=TRUE)[1]
```

```
## Liverpool
```

```
## 8
```

```
sort(table(survey$bev),decreasing=TRUE)[1]
```

```
## coffee
```

```
## 19
```

28. Using a logical subscript, print out the distance from Liverpool city centre of everyone who listed Liverpool as their top football club.

```
survey$distance[survey$football_club=='Liverpool' & !is.na(survey$football_club)]
```

```
## [1] 0.0 10.0 4.0 1.0 8.0 0.6 1.5 3.0
```

29. Print out the distances of everyone who listed a football club other than Liverpool as their top football club. (This should *not* include those with no preference / no answer!)

```
survey$distance[survey$football_club!='Liverpool' & survey$football_club!='none' &
!is.na(survey$football_club)]
```

```
## [1] 5.0 4.0 0.7 2.0 3.0 0.8 0.7 1.0 1.0 0.0 0.5 4.0 15.0 1.0
```

30. Calculate the average height of women in the survey, and the average height of men in the survey.

```
mean(survey$height[survey$gender=='female'],na.rm=TRUE)
```

```
## [1] 64.41667
```

```
mean(survey$height[survey$gender=='male'],na.rm=TRUE)
```

```
## [1] 69.88889
```

```
# - or -
```

```
tapply(survey$height,survey$gender,mean,na.rm=TRUE)
```

```
##      female      male
```

```
## 64.41667 69.88889
```

31. Calculate the mean heights of everyone in the survey broken down by their preferred beverage (i.e.: the mean heights of coffee-drinkers, tea-drinkers, etc.) Note: You can do this in a single line in R if you use `tapply`.

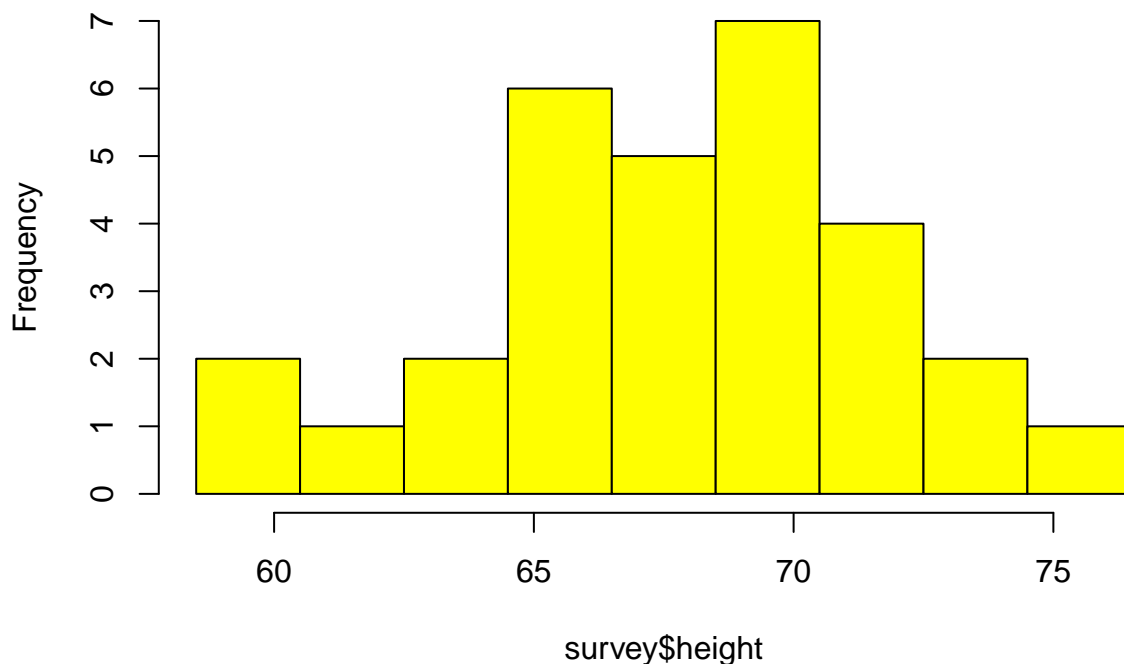
```
tapply(survey$height,survey$bev,mean,na.rm=TRUE)
```

```
##      coffee      cola energy drink      none      tea
## 68.11111 69.33333 68.00000 67.00000 63.60000
```

32. Produce a histogram of (all) student heights. Note: be careful about the default break points!

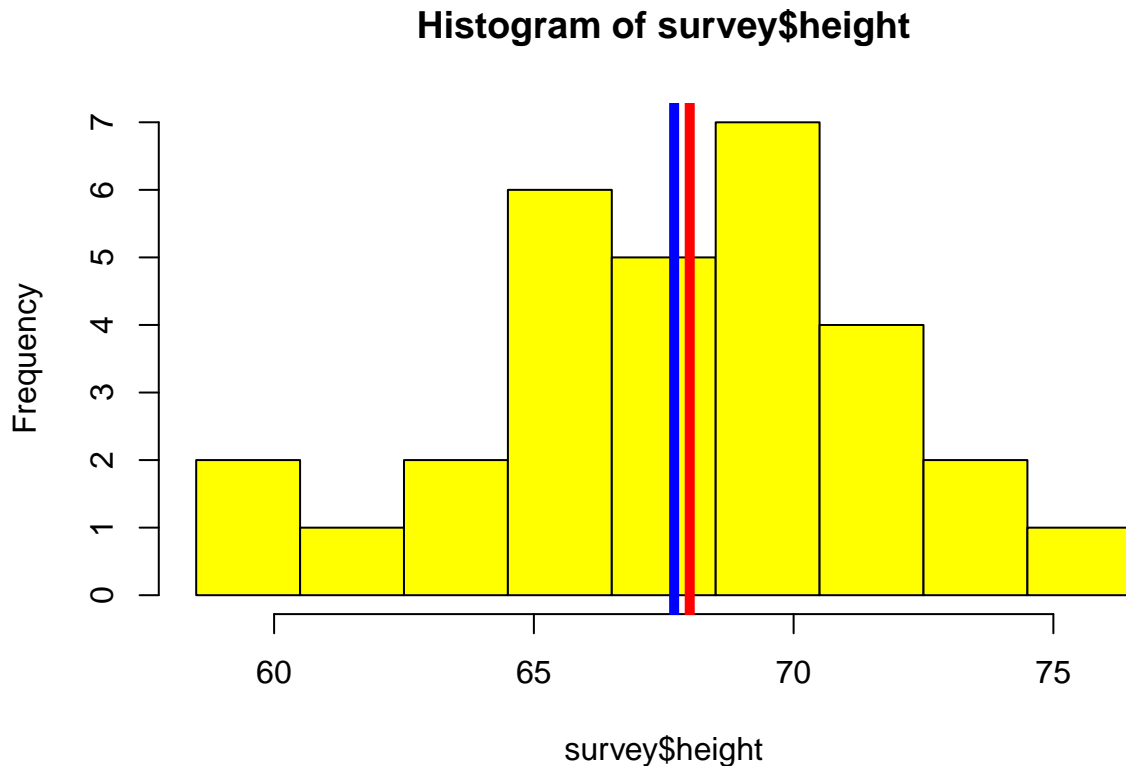
```
hist(survey$height,breaks=seq(58.5,76.5,2),col='yellow')
```

Histogram of survey\$height



33. Overplot (plot on top of the previous plot) two vertical lines of different colours: one at the MEAN value, and one at the MEDIAN value. (Note: in Rmarkdown you'll need to repeat the plotting command from #32 in your code chunk for #33.)

```
hist(survey$height,breaks=seq(58.5,76.5,2),col='yellow')
abline(v=mean(survey$height,na.rm=TRUE),lwd=5,col='blue')
abline(v=median(survey$height,na.rm=TRUE),lwd=5,col='red')
```



34. Suppose that one of the students who did not specify a football club preference decides to support a team. Choose a student that did not specify a football club (or specified 'none') and update the dataframe by assigning a football team of your choice to that student. Then print out an updated frequency table containing the number of students supporting each football club.

```
surveymod = survey
which_noclub = which(surveymod$football_club=='none')
surveymod$football_club[which_noclub[1]] = 'Liverpool'
table(surveymod$football_club)
```

```
##
##      Arsenal      Barcelona      Chelsea      Leeds      Liverpool
##           1           1           4           1           9
## Manchester City      none      none      PSG      Real Madrid
##           1           1           2           1           2
##      Wisla Krakow
##           1
```

Larger Data Frames and Normality

The file `abalone.csv` (available on Canvas) contains data from a sample of blacklip abalones (*Haliotis rubra*) gathered from the Bass Straits of Tasmania in the early 1990s. The measurements include various dimensions of the shell (length, diameter, height) and weight after various levels of treatment. “Rings” gives the number of ring layers inside the shell and is a proxy for age since about one ring layer is added per year. Dimensions are in metres and weights are in kilograms.

35. Download the file `abalone.csv` to your local computer. Load it into R as a data frame, and print out the first few lines using `head`.

```
abalone = read.csv('abalone.csv',as.is=FALSE)
head(abalone)

##   sex length diameter height weight shucked.weight viscera.weight shell.weight
## 1  M  0.455    0.365  0.095 0.5140          0.2245          0.1010          0.150
## 2  M  0.350    0.265  0.090 0.2255          0.0995          0.0485          0.070
## 3  F  0.530    0.420  0.135 0.6770          0.2565          0.1415          0.210
## 4  M  0.440    0.365  0.125 0.5160          0.2155          0.1140          0.155
## 5  I  0.330    0.255  0.080 0.2050          0.0895          0.0395          0.055
## 6  I  0.425    0.300  0.095 0.3515          0.1410          0.0775          0.120
##   rings
## 1     15
## 2      7
## 3      9
## 4     10
## 5      7
## 6      8
```

36. How many different variables (columns) are in the table, what are their names, and what sort of data do they contain? (Continuous, categorical, others)? How many measurements are there in the sample? (Use the `str` command.)

```
str(abalone)

## 'data.frame':   4177 obs. of  9 variables:
##  $ sex          : Factor w/ 3 levels "F","I","M": 3 3 1 3 2 2 1 1 3 1 ...
##  $ length       : 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 ...
##  $ weight       : num  0.514 0.226 0.677 0.516 0.205 ...
##  $ shucked.weight: num  0.2245 0.0995 0.2565 0.2155 0.0895 ...
##  $ viscera.weight: num  0.101 0.0485 0.1415 0.114 0.0395 ...
##  $ shell.weight  : 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 ...
```

37. Produce a summary table of the abalone dataset using `summary`. Make sure you understand what all the numbers mean.

```
summary(abalone)

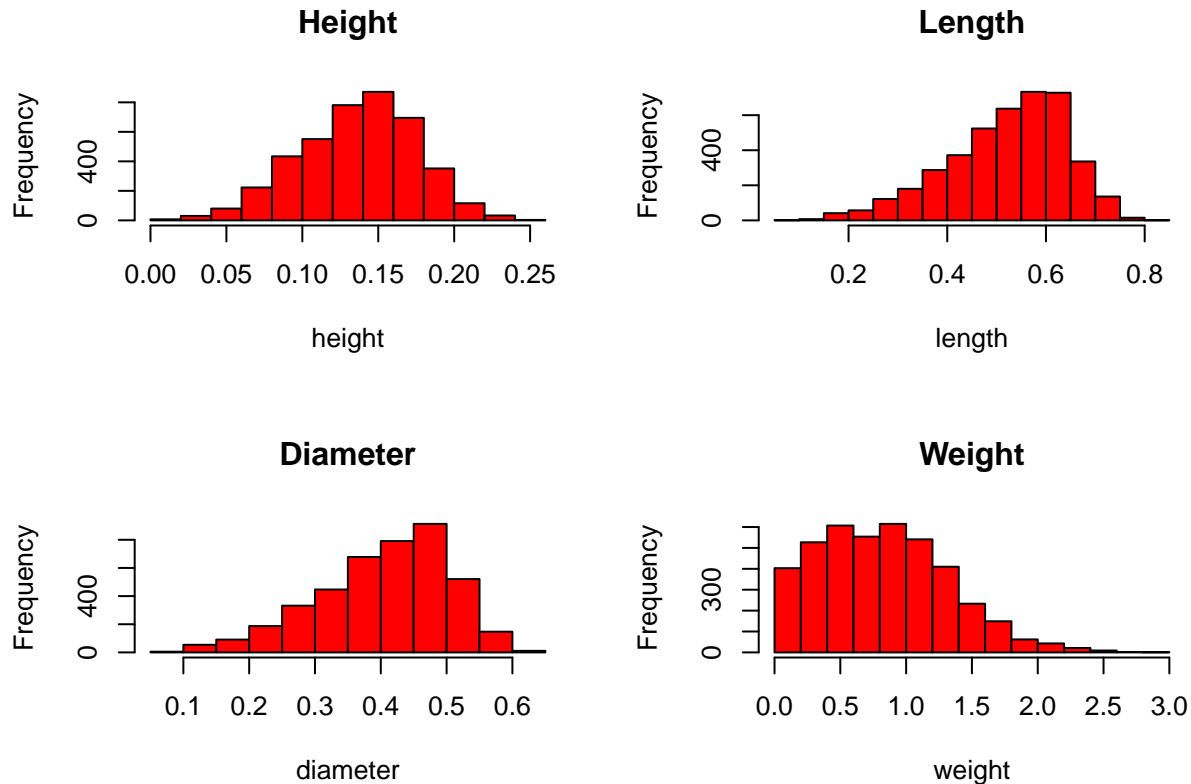
##   sex          length      diameter      height      weight
## F:1307   Min.    :0.075   Min.    :0.0550   Min.    :0.0000   Min.    :0.0020
## I:1342   1st Qu.:0.450   1st Qu.:0.3500   1st Qu.:0.1150   1st Qu.:0.4415
## M:1528   Median :0.545   Median :0.4250   Median :0.1400   Median :0.7995
##          Mean   :0.524   Mean   :0.4079   Mean   :0.1392   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.   :0.2500   Max.   :2.8255
##                                     NA's    :1
## shucked.weight viscera.weight shell.weight 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
##
```

38. Produce a four-panel set of plots (use `par(mfrow)` to set this up) containing histograms of the following four quantities: height, length, diameter, weight. Make the axis labels and titles reader-friendly and use colours to fill the histogram bars.

```
par(mfrow=c(2,2))
hist(abalone$height,col='red',main='Height',xlab='height')
hist(abalone$length,col='red',main='Length',xlab='length')
hist(abalone$diameter,col='red',main='Diameter',xlab='diameter')
hist(abalone$weight,col='red',main='Weight',xlab='weight')
```



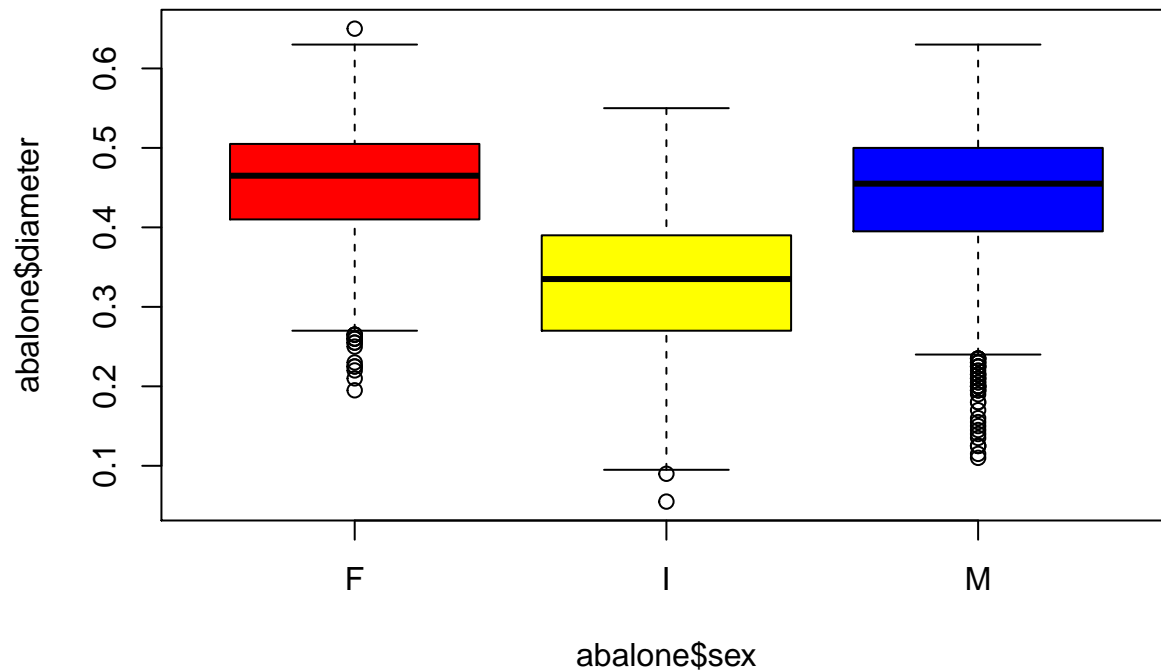
```
par(mfrow=c(1,1))
```

39. Do any of these quantities appear to be normally distributed?

Height might be close to normal, but length, diameter, and weight are clearly not normal.

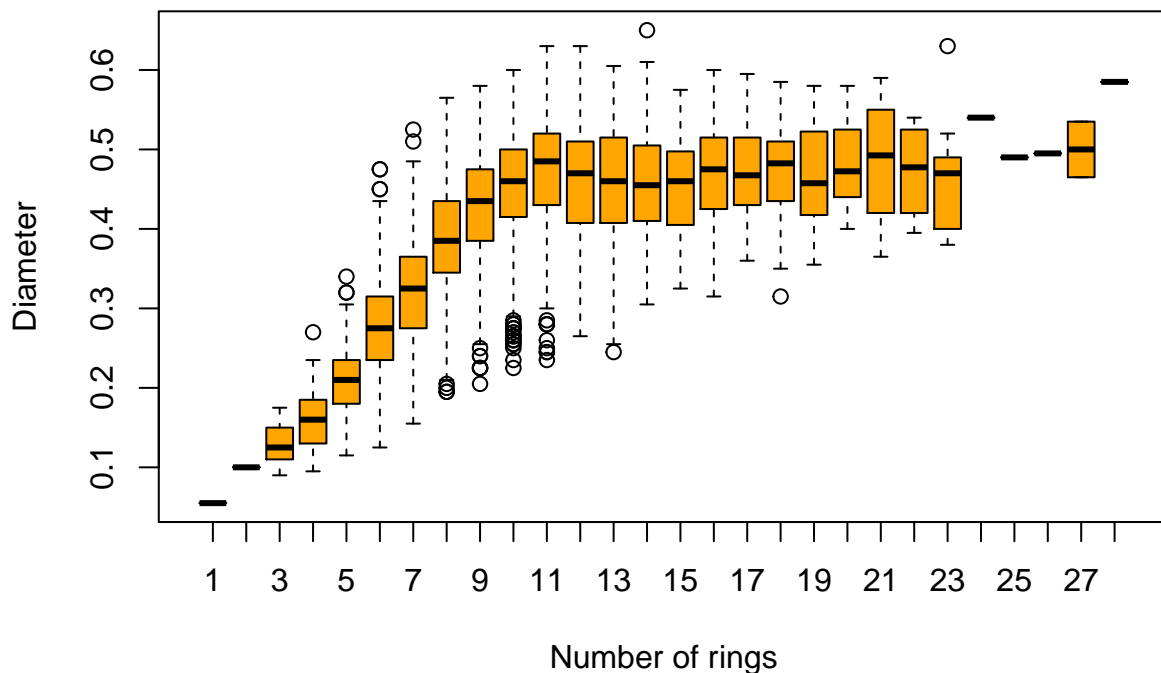
40. Produce a (one-panel) box-and-whisker plot summarising the distribution of abalone diameter for the three sex categories (F, I, or M). Be sure to label the axes. Do the distributions of females and males seem obviously different? What about for “I”, meaning indeterminate/unknown sex?

```
boxplot(abalone$diameter~abalone$sex,col=c('red','yellow','blue'))
```



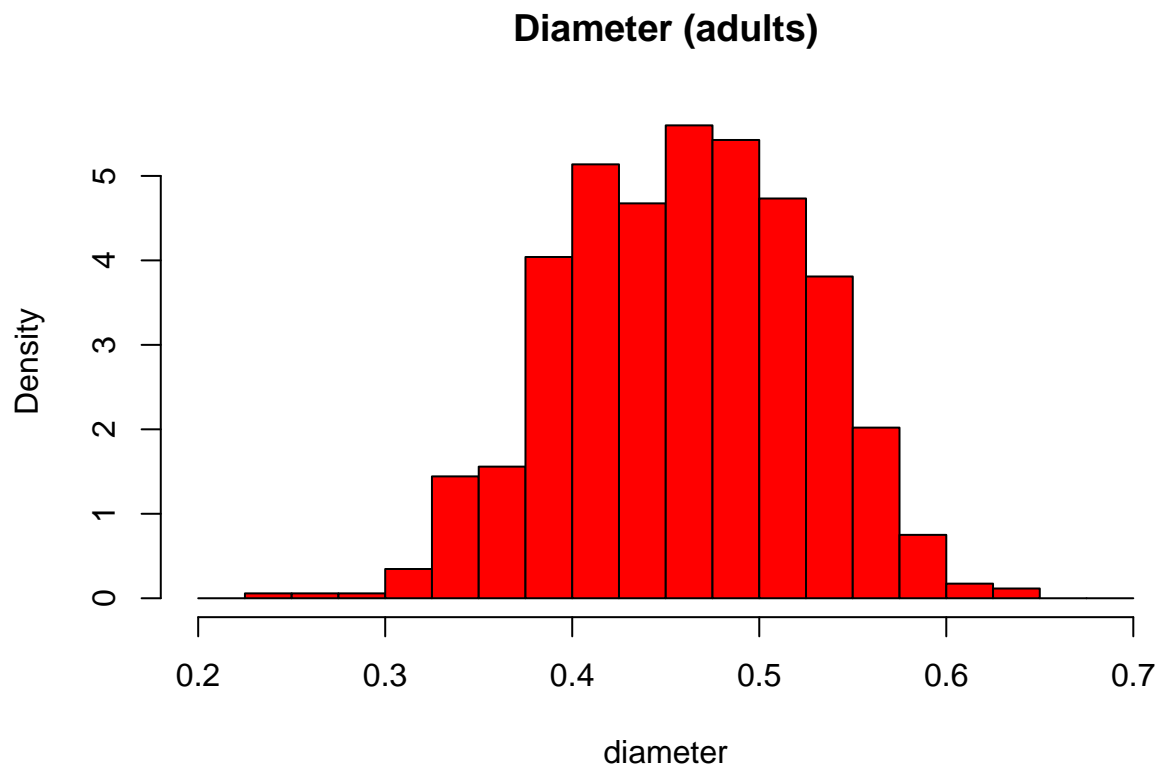
41. Produce a box-and-whisker plot showing the distribution of abalone diameter as a function of number of rings (a measurement proxy for age in years.)

```
boxplot(abalone$diameter~abalone$rings,col='orange',xlab='Number of rings',ylab='Diameter')
```



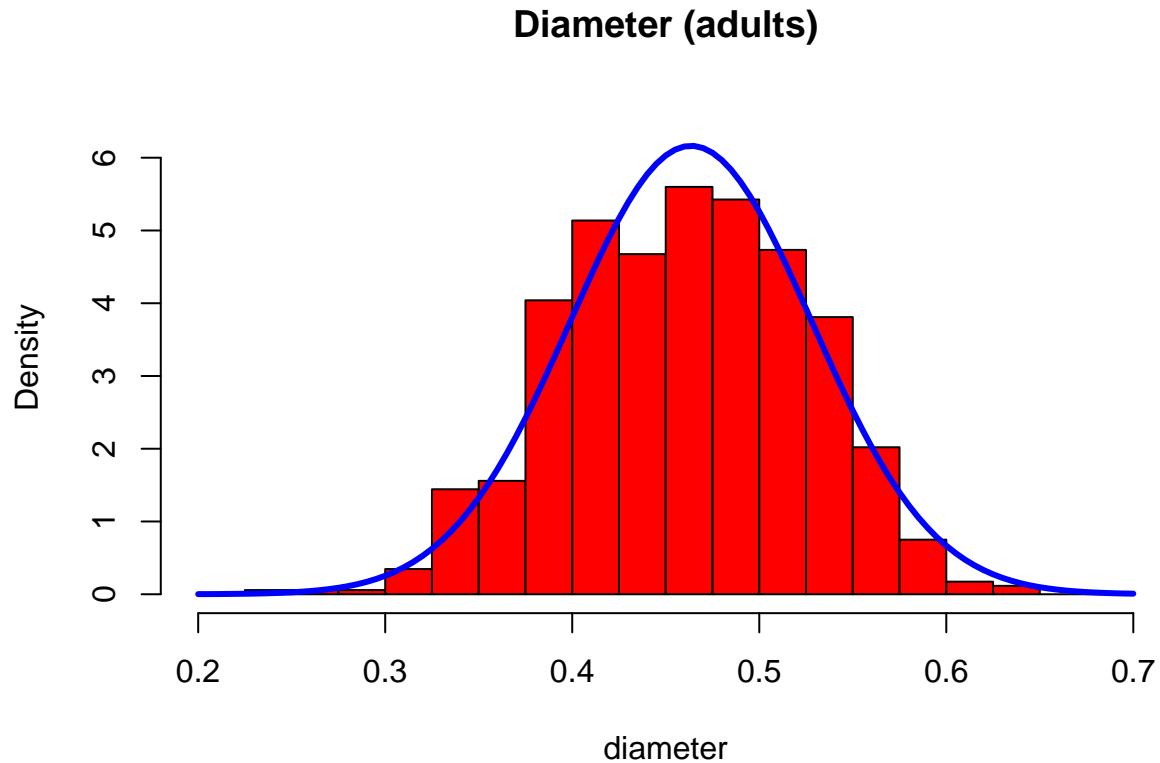
42. Create a new dataframe that contains a subset of the original data frame: specifically, it includes only those rows for which nrings is greater than or equal to 13. Then, produce a histogram of the diameters of this subset. (Set `freq=FALSE` to plot probability density instead of frequency, and you may want to specify the breaks to get better resolution.) Do these look normally distributed?

```
adults = abalone[abalone$rings >= 13,]
hist(adults$diameter,col='red',main='Diameter (adults)',xlab='diameter',breaks=seq(0.2,0.7,0.025),freq=
```



43. Check your assessment of normality by over-plotting a normal distribution on top of the histogram above, using the appropriate μ and σ parameters to describe the population. (Start by defining a plotting variable along the x-axis using `seq`, and then calculate the normal PDF as a function of this variable using the function `dnorm`.)

```
hist(adults$diameter,col='red',main='Diameter (adults)',xlab='diameter',
     breaks=seq(0.2,0.7,0.025),freq=FALSE,ylim=c(0,6.5))
d = seq(0.2,0.7,0.005)
m = mean(adults$diameter)
s = sd(adults$diameter)
lines(d,dnorm(d,mean=m,sd=s),col='blue',lwd=3)
```



44. Explain why the central limit theorem did not seem to hold for the entire data set, even though it did for a subset.

For juvenile abalone the diameter is a function of largely a single parameter: their age. As a result the distribution primarily reflects the age distribution (which is not normal) rather than a combination of many factors. If only adult abalone are considered, growth has stopped and their final diameter largely a product of complex, additive genetic and environmental factors, so a normal distribution is produced.