
Using R: Frequency Distributions, Histograms, Scatterplots, & Line Graphs

This document describes how to accomplish the following tasks.

- **Making a Frequency Table** `table(data)` 2
Here we create a frequency table from raw data imported from a .CSV file. We also see how to append a relative and cumulative frequency table to the original frequency table.
- **Making Histograms** `hist(data)` 4
Here use the `hist` command to make a fast and dirty histogram and demonstrate how to add some bells and whistles.
- **Making Scatterplots** `plot(x-data, y-data)` 5
Using the `plot` command.
- **Making Line Graphs and Time Series Plot** `plot(time-data, y-data)` 6
Using the `plot` command when the x -values are dates.
Here we use the `as.Date` command to correctly read dates.

- **Entering Data**
 - **Making a Data List (vector)** `list-name <- c(#, #, #, ..., #)` 7
 - **Making a Table (matrix)** `table-name <- matrix(....)` 7
 - **Importing Data from a CSV File** 8
`table-name <- read.csv(file="file-name.csv", ...)`
or
`table-name <- read.table("File-Name.csv",header=TRUE,sep=",")`

- **Frequency Tables in R:** In the textbook, we took 42 test scores for male students and put the results into a frequency table.

Scores on Test #2 - Males						
42 Scores: Average = 73.5						
84	88	76	44	80	83	51
93	69	78	49	55	78	93
64	84	54	92	96	72	97
37	97	67	83	93	95	67
72	67	86	76	80	58	62
69	64	82	48	54	80	69

Raw Data
→ becomes →
Organized

Males	
Scores	Frequency
30 - 39	1
40 - 49	3
50 - 59	5
60 - 69	9
70 - 79	6
80 - 89	10
90 - 99	8

Then we created a relative and cumulative frequency table from this.

Frequency Distribution: Males		Relative Frequency Distribution: Males		Cumulative Frequency Distribution: Males	
Scores	Frequency	Scores	Relative Frequency	Scores	Cumulative Frequency
30 - 39	1	30 - 39	2.4%	less than 40	1
40 - 49	3	40 - 49	7.1%	less than 50	4
50 - 59	5	50 - 59	11.9%	less than 60	9
60 - 69	9	60 - 69	21.4%	less than 70	18
70 - 79	6	70 - 79	14.3%	less than 80	24
80 - 89	10	80 - 89	23.8%	less than 90	34
90 - 99	8	90 - 99	19.0%	less than 100	42

Here we see how to do these tasks with R. We'll start by importing the data into R. Suppose the data is in an Excel file saved as a .CSV file named `Excel-Data.csv` that looks like

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	Males	Females												
2		84	59											
3		93	58											
4		64	64											
5		37	96											
6		72	79											

This is imported into R using the `read.table` command.

```
> male_female_data <- read.table("Excel-Data.csv", header=TRUE, sep=",")
> male_scores <- male_female_data$Males
```

Here we have R create a frequency table and then append a relative and cumulative table to it. Everything in red is typed by the user. Everything in blue is output to the console.

```

> male_female_data <- read.table("Excel-Data.csv",header=TRUE,sep=",")
> male_scores <- male_female_data$Males
> bins <- seq(29.5,99.5,by=10) # creates a list of class boundaries
> Scores <- cut(male_scores, bins) # group the data into bins
> ### below produces the original table in a bad format
> table(Scores)
Scores
(29.5,39.5] (39.5,49.5] (49.5,59.5] (59.5,69.5] (69.5,79.5] (79.5,89.5]
      1          3          5          9          6         10
(89.5,99.5]
      8
> ### below produces the same table in a better format
> transform(table(Scores))
      Scores Freq
1 (29.5,39.5]    1
2 (39.5,49.5]    3
3 (49.5,59.5]    5
4 (59.5,69.5]    9
5 (69.5,79.5]    6
6 (79.5,89.5]   10
7 (89.5,99.5]    8
> #### below appends a relative and cumulative table.
> transform(freq_table,Rel_Freq=prop.table(Freq),Cum_Freq=cumsum(Freq))
      Scores Freq  Rel_Freq Cum_Freq
1 (29.5,39.5]    1 0.02380952     1
2 (39.5,49.5]    3 0.07142857     4
3 (49.5,59.5]    5 0.11904762     9
4 (59.5,69.5]    9 0.21428571    18
5 (69.5,79.5]    6 0.14285714    24
6 (79.5,89.5]   10 0.23809524    34
7 (89.5,99.5]    8 0.19047619    42
> |

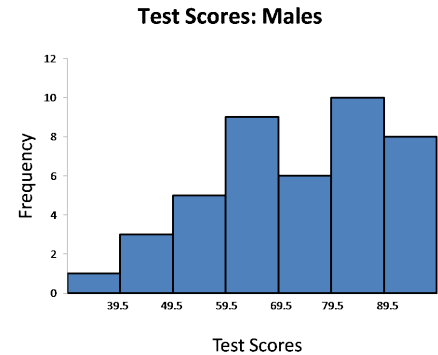
```

- The classes are defined by creating a list of class boundaries. You can create this list by hand or
`> bins <- seq(29.5,99.5,by=10)`
 The `seq` function creates a list by starting at 29.5 and increasing by 10 until it hits 99.5. Now `bins = [29.5, 39.5, 49.5, 59.5, 69.5, 79.5, 89.5, 99.5]`.
- The `cut` function organizes the data into the appropriate bins.
- The `transform` function puts the tables into column format - which is nice.

- **Histograms in R:** In the text, we created a histogram from the raw data.

Scores on Test #2 - Males						
42 Scores: Average = 73.5						
84	88	76	44	80	83	51
93	69	78	49	55	78	93
64	84	54	92	96	72	97
37	97	67	83	93	95	67
72	67	86	76	80	58	62
69	64	82	48	54	80	69

Raw Data
→ becomes →
Histogram

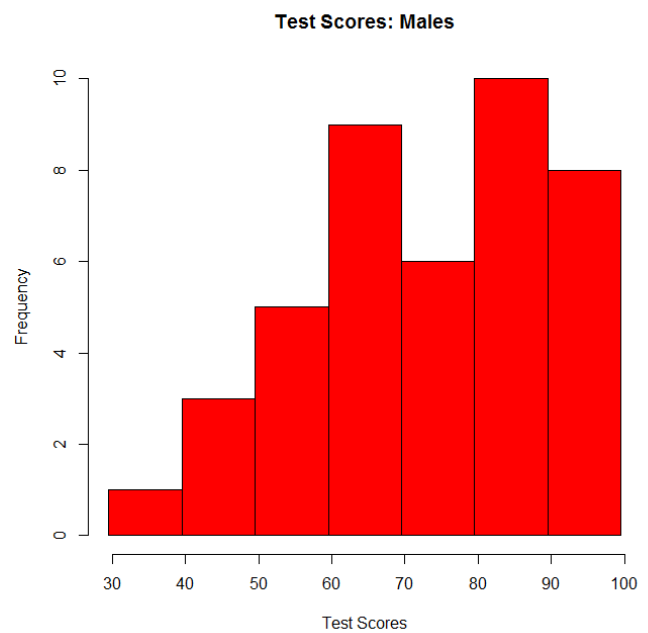
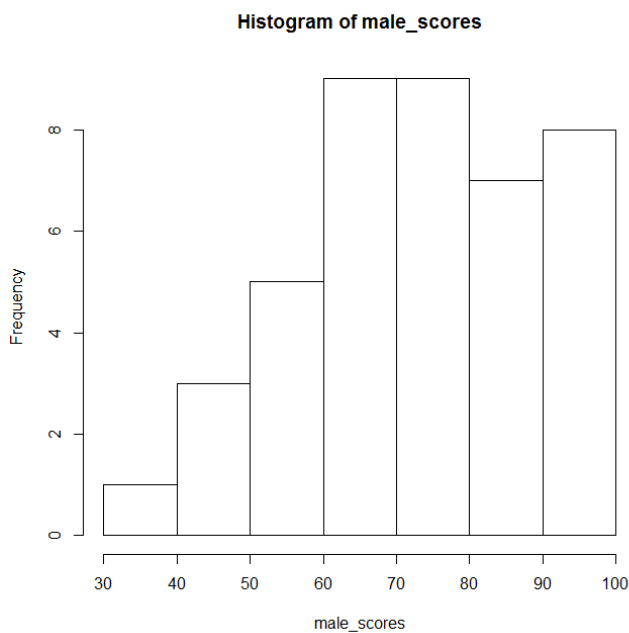


Here, we'll let R create the histogram using the `hist` command. You can define your own classes by creating a list of class boundaries and using the `breaks =` command. You can also add a title (`main =`), a label (`xlab =`), and color (`col =`).

```

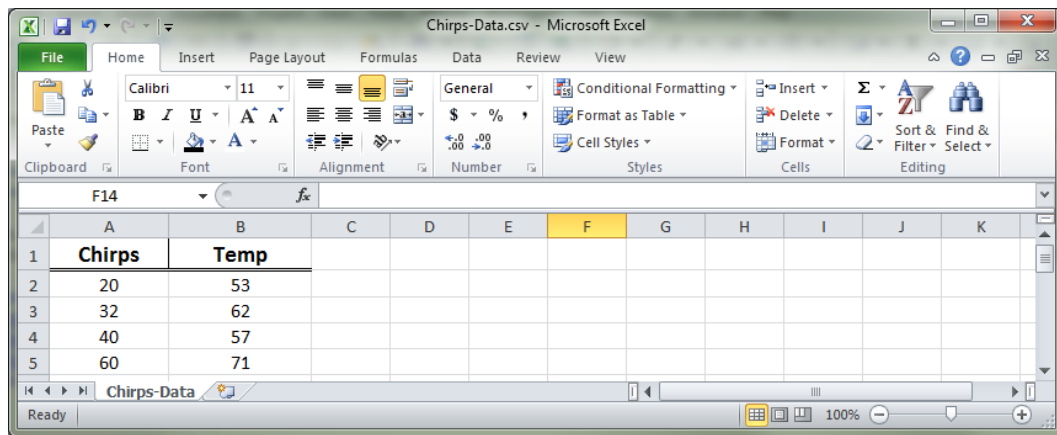
RGui (32-bit) - [R Console]
File Edit View Misc Packages Windows Help
[Icons]
> male_female_data <- read.table("Excel-Data.csv",header=TRUE,sep=",")
> male_scores <- male_female_data$Males
>
> ## fast and dirty (below left)
> hist(male_scores) # creates the histogram below left
>
> ## Add some bells and whistles (below right)
> bins <- seq(29.5,99.5,by=10) #creates a list of class boundaries
> hist(male_scores,breaks=bins,main="Test Scores: Males", xlab = "Test Scores",col="red")
> |

```



• Scatterplots in R:

Suppose we have data for cricket chirps per minute and temperature in degrees Fahrenheit in an Excel file saved in .CSV format that looks like



	A	B	C	D	E	F	G	H	I	J	K
1	Chirps	Temp									
2	20	53									
3	32	62									
4	40	57									
5	60	71									

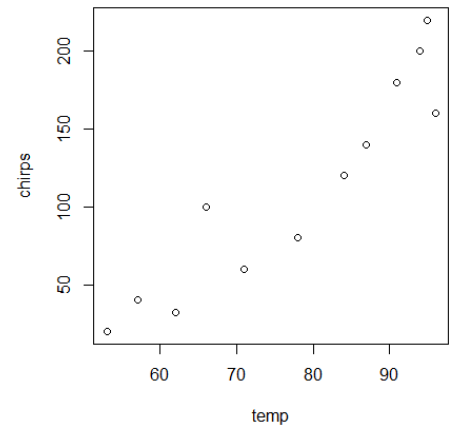
We have R create a scatterplot with the `plot(x,y)` command and put in the line of best fit with the `abline` command.

Simple scatterplot:

```

RGui (32-bit) - [R Console]
File Edit View Misc Packages Windows Help
> chirps_data <- read.table("Chirps-Data.csv",header=TRUE,sep=",")
> temp <- chirps_data$Temp
> chirps <- chirps_data$Chirps
> plot(temp,chirps)
> |

```

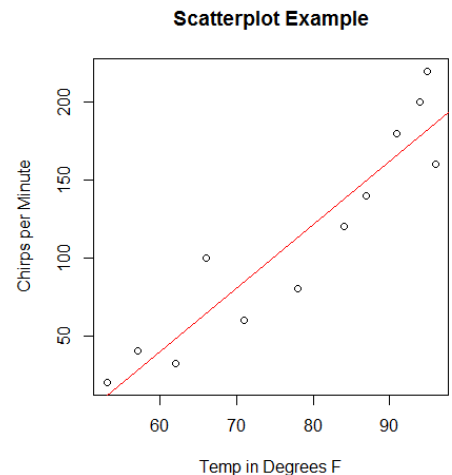


Some customization and a line of best fit:

```

RGui (32-bit) - [R Console]
File Edit View Misc Packages Windows Help
> ##### customize and add a line of best fit
> plot(temp,chirps, main="Scatterplot Example",
+       xlab="Temp in Degrees F", ylab="Chirps per Minute")
> abline(lm(chirps~temp), col="red") #regression line (y~x)
> |

```



• Line Graphs and Time Series Graphs in R:

A line graph is just a scatterplot where the points are connected moving left to right. Getting the points connected is done using the `type` command.

```
> plot(x-data,y-data,type="o")
```

This now puts a small circle at each point and then connects the points with a line. Any `type` will result in the points being connected with a line.

A time series graph is a line graph where the x -axis represents time. If you are doing a time series with clock time (seconds, minutes, hours), you just create a line graph with the appropriate time units on the x -axis. The tricky part is handling dates. Suppose you have your dates in an Excel spreadsheet saved as a CSV file like the one below. In this example we look at the unemployment rates from July 2008 - July 2009 (the first 12 months after the financial collapse of 2007).

	A	B	C	D	E	F	G	H	I	J	K	L
1	Date	Unemployment_Rate										
2	7/1/2008	5.8										
3	8/1/2008	6.1										
4	9/1/2008	6.2										
5	10/1/2008	6.6										

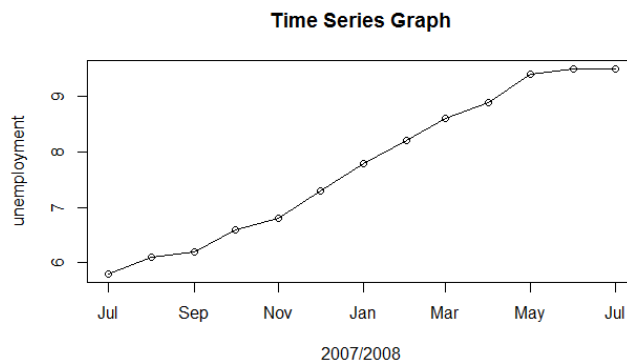
The trick is to tell R how to read the date data by

```
> data <- read.table("Unemployment-Data.csv",header=TRUE,sep=",")
```

```
> month <- as.Date(data$Date, "%m/%d/%Y")
```

We have R create a time series graph with the `plot` command.

```
> data <- read.table("Unemployment-Data.csv",header=TRUE,sep=",")
> month <- as.Date(data$Date, "%m/%d/%Y")
> unemployment <- data$Unemployment_Rate
> plot(month,unemployment,type="o",
+       xlab="2007/2008",main="Time Series Graph")
> |
```



- **Making a Data List (vector):** When entering data you must use the syntax below.

```
> Mon <- c(68, 84, 93, 68, 70)
```

- Here, `Mon` is the name of the list (you get to choose this name).
 - Then, you have the 'less than sign' from the keyboard followed by the dash ('-') symbol.
 - All lists must be enclosed with `c(#, #, #, ...)`. I think the `c` stands for 'combine'.
 - Don't forget to separate the data values with commas.
- **Making a Table (matrix):** Suppose you want to make a table of 5-numbers for each of the first three days of the work-week.

	Monday	Tuesday	Wednesday
	68	59	55
	84	72	66
	93	78	77
	68	91	88
	70	90	99

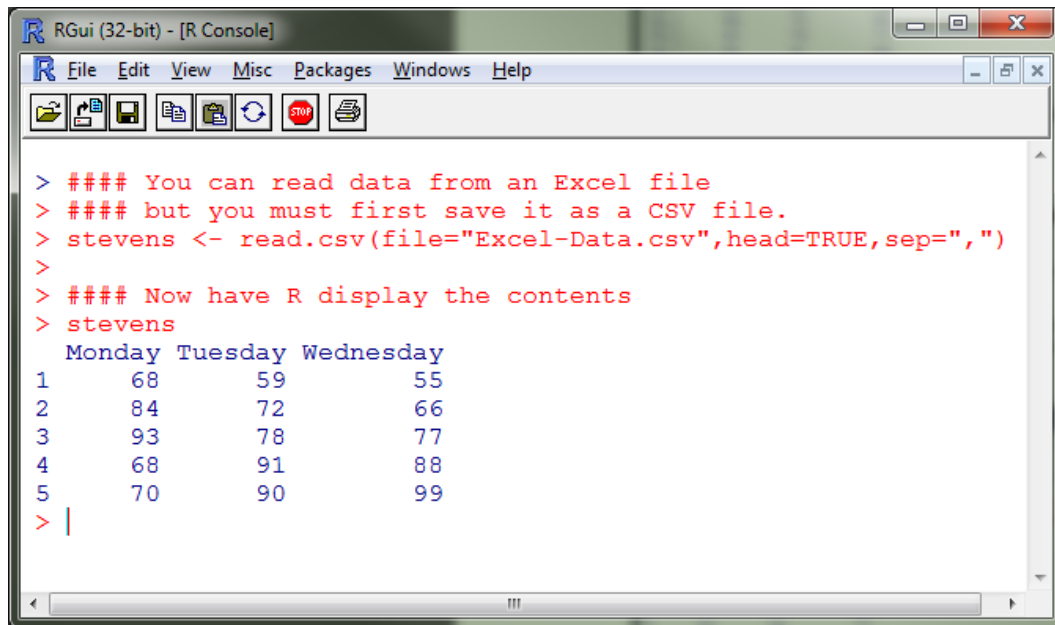
You can combine lists like the one described above into a table using the `matrix` command. Setting this up takes a little effort and the usage is depicted in the screenshot on the next page. **Everything in red is typed by the user.** **Everything in blue is output to the console.**

```

RGui (32-bit) - [R Console]
File Edit View Misc Packages Windows Help
[Icons]
> Mon <- c(68, 84, 93, 68, 70)
> Tue <- c(59, 72, 78, 91, 90)
> Wed <- c(55, 66, 77, 88, 99)
>
>
> ##### Put these into a table where each day gets its own column:
> cells      <- c(Mon, Tue, Wed)
> column_names <- c("Monday", "Tuesday", "Wednesday")
> Table      <- matrix(cells, nrow=5, ncol=3, byrow=FALSE)
> colnames(Table) = column_names;
>
> ##### Now have R display the contents
> Mon
[1] 68 84 93 68 70
> Tue
[1] 59 72 78 91 90
> Wed
[1] 55 66 77 88 99
> Table
      Monday Tuesday Wednesday
[1,]     68      59         55
[2,]     84      72         66
[3,]     93      78         77
[4,]     68      91         88
[5,]     70      90         99
> |

```

- **Importing Data:** Since all of the other software packages will easily convert a data file into a CSV file, we will use this format to read the data into R. The screenshot below depicts how to read such a file and display the contents. The `head=TRUE` means the first row contains column headings (not data).



```
RGui (32-bit) - [R Console]
File Edit View Misc Packages Windows Help
[Icons]

> #### You can read data from an Excel file
> #### but you must first save it as a CSV file.
> stevens <- read.csv(file="Excel-Data.csv",head=TRUE,sep=",")
>
> #### Now have R display the contents
> stevens
  Monday Tuesday Wednesday
1     68      59         55
2     84      72         66
3     93      78         77
4     68      91         88
5     70      90         99
> |
```