

STAT167 Lab #1 - Spring 2025

Ethan Choi

2025/4/4

Contents

Discussion/Lab #1 instructions	2
RMarkdown Basics	2
Test Header 1	3
Test Header 2	3
Test Header 3	3
Packages	3
Install necessary packages	3
Adding R	3
R Code Chunks	3
Inline R	4
Chunk Options	4
Adding Math with LaTeX	4
Displaystyle LaTeX	4
Inline LaTeX	4
Lecture Review	5
Generate sample data	5
Exercise #1	6
Calculate summary statistics	7
Exercise #2	7
Exercise #3	7
Plot single data set	8
Plot a theoretical standard normal curve	9
Exercise #4	11

Acknowledgment: Part of this R Markdown template is adapted from David Dalpiaz (UIUC) and Cosma Shalizi (CMU).

Discussion/Lab #1 instructions

This week, you will first learn the basic R Markdown syntax.

- First, download the `rmd` file from Canvas.
- Open this `rmd` file in RStudio and click **Knit -> Knit to PDF** to render it to PDF format. You need to have **LaTeX** installed on the computer to render it to PDF format. If not, you can also render it to HTML format.
- Read this `rmd` file and the rendered `pdf/html` file side-by-side, to see how this document was generated!
- Be sure to play with this document! Change it. Break it. Fix it. The best way to learn R Markdown (or really almost anything) is to try, fail, then find out what you did wrong.

Next, you will review some example code from this week's lectures.

- Read over the code and the output. If you have any questions about certain functions or parameters, it is the time to ask!
- There are some exercises through out this document. Replace **INSERT_YOUR_ANSWER** with your own answers. Knit the file, and check your results.

Please comment your R code thoroughly, and follow the R coding style guideline (<https://google.github.io/styleguide/Rguide.xml>). Partial credit may be deducted for insufficient commenting or poor coding styles.

Lab submission guideline

- After you completed all exercises, save your file to `FirstnameLastname-SID-lab1.rmd` and save the rendered pdf file to `FirstnameLastname-SID-lab1.pdf`. If you can not knit it to pdf, knit it to html first and then print/save it to pdf format.
- Submit **BOTH your source rmd file and the knitted pdf file to GradeScope**. Do NOT create a zip file.
- You can submit multiple times, you last submission will be graded.

RMarkdown Basics

RMarkdown at its core is a combination of R and Markdown used to generate reproducible reports for data analyses.

Markdown and R are mixed together in a `.rmd` file, which can then be rendered into a number of formats including `.html`, `.pdf`, and `.docx`. There is a strong preference for using `.html` in this course.

Have a look at this `.rmd` to see how this document was generated! It should be read alongside the rendered `.html` to best understand how everything works. You can also modifying the `.rmd` along the way, and see what effects your modifications have.

Formatting text is easy. **This is bold.** *This is italics.* This text appears as monospaced.

Test Header 1

Test Header 2

Test Header 3

- Unordered list element 1.
- Unordered list element 2.
- Unordered list element 3.

1. Ordered list element 1.
2. Ordered list element 2.
3. Ordered list element 3.

Packages

Packages are key to using R. The community generated packages are a large part of R's success, and it is extremely rare to perform an analysis without using at least some packages. Once installed, packages must be loaded before they are used.

Install necessary packages

```
# If you can knit the html file successfully, do not change this R code chunk  
  
# Otherwise, uncomment the two lines and install the packages.  
# You only need to run the install.packages() commands once.  
# Then you can comment them out by adding `#` back at the beginning of each line  
  
#install.packages("rmarkdown", repos="http://cran.rstudio.com/")  
#install.packages("yaml", repos="http://cran.rstudio.com/")
```

Note that `rmarkdown` is actually a package in R! If R never prompts you to install `rmarkdown` and its associated packages when first creating an RMarkdown document, use the above command to install them manually.

Adding R

So far we have only used Markdown to create html. This is useful by itself, but the real power of RMarkdown comes when we add R. There are two ways we can do this. We can use R code chunks, or run R inline.

R Code Chunks

The following is an example of an R code chunk

```
# generate random normals  
set.seed(167) # feel free to change 167 to your lucky number  
x <- rnorm(100) # use <- for assignment instead of =
```

There is a lot going on here. In the `.rmd` file, notice the syntax that creates and ends the chunk. Also note that `example_chunk` is the chunk name. Everything between the start and end syntax must be valid **R** code. Chunk names are not necessary, but can become useful as your documents grow in size.

Inline R

R can also be ran in the middle of exposition. For example, the mean of the data we generated is -0.094945.

Chunk Options

There are many chunk options. Here we first introduce two options which are frequently used: `eval` and `echo`.

```
?log
x
```

Using `eval = FALSE` the above chunk displays the code, but it is not run. The `?` code pulls up documentation of a function. The `x` code prints all values inside the vector `x`. With `eval = FALSE` (or simply `eval = F`), there is no output displayed.

```
## [1] "Hello World!"
```

Above, we see output, but no code! This is done using `echo = FALSE` (or simply, `echo=F`), which is often useful.

Adding Math with LaTeX

Another benefit of RMarkdown is the ability to add Latex for mathematics typesetting. Like **R** code, there are two ways we can include Latex; `displaystyle` and `inline`.

Note that use of **LaTeX** is somewhat dependent on the resulting file format. For example, it cannot be used at all with `.docx`. To use it with `.pdf` you must have LaTeX installed on your machine.

With `.html` the **LaTeX** is not actually rendered during knitting, but actually rendered in your browser using MathJax.

Displaystyle LaTeX

Displaystyle is used for larger equations which appear centered on their own line.

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

Inline LaTeX

We could mix LaTeX commands in the middle of exposition, for example: $t = 2$. We could actually mix **R** with Latex as well! For example: $\bar{x} = -0.094945$.

Lecture Review

Generate sample data

Recall that in the lecture, we generated 100 random normals, assigned the values into a vector `x`, then printed out some values in vector `x`.

```
# generate random normals
set.seed(167) # feel free to change 167 to your lucky number
x <- rnorm(100) # use <- for assignment instead of =
# print out the data or part of it
x
```

```
## [1] -0.99307913 -0.17111961 -0.10200147 -0.14990885 -0.21135557 -0.60078127
## [7] 0.11739337 -0.21578390 -0.75237199 1.55294726 0.77671331 0.32967895
## [13] -1.89405513 -1.99349413 -0.22007207 0.19850126 0.56371430 -0.88079442
## [19] -0.64174738 -0.84737192 1.83233164 -0.53413321 0.30973709 0.01403343
## [25] -1.09327949 -0.07513613 1.52135816 0.53832708 -0.26987651 0.26677743
## [31] -0.70915890 -1.27608985 0.52580565 0.71607236 -0.60290030 0.14002073
## [37] -1.99917225 0.42774676 -0.08996853 -0.23563275 -0.41061735 0.15146705
## [43] -2.26942347 -0.93500192 -1.41073822 -2.32958937 -0.15649706 -0.94767522
## [49] 0.73027385 0.41461206 0.65270923 0.64319947 -0.08089847 0.27330967
## [55] -0.08968567 -0.37876327 -1.15709419 0.58358139 -0.07484820 0.85993573
## [61] 0.70481069 0.66931003 0.38240296 0.01418130 1.05815293 -0.97596451
## [67] -0.01815480 -1.60150608 -0.29764517 -0.97331417 -1.53441965 2.09839585
## [73] 0.95002228 0.09173245 1.54627570 -1.42063598 0.73546857 0.12136776
## [79] 1.03977273 -1.38440186 0.05776784 0.17496418 -0.19319262 0.06309704
## [85] 0.52483624 0.12168457 1.60986947 0.96609625 -0.89391965 -0.89837740
## [91] -1.56742899 1.73821583 -1.21415287 1.44463281 -1.61656091 0.40198802
## [97] 0.15775583 0.57040588 0.07979506 0.43204664
```

```
1:10
```

```
## [1] 1 2 3 4 5 6 7 8 9 10
```

```
10:1
```

```
## [1] 10 9 8 7 6 5 4 3 2 1
```

```
seq(from=1, to=10, by=2)
```

```
## [1] 1 3 5 7 9
```

```
seq(10, 1)
```

```
## [1] 10 9 8 7 6 5 4 3 2 1
```

```
x[1:10]
```

```
## [1] -0.9930791 -0.1711196 -0.1020015 -0.1499089 -0.2113556 -0.6007813
## [7] 0.1173934 -0.2157839 -0.7523720 1.5529473
```

```
y <- c(2,4,7)
y
```

```
## [1] 2 4 7
```

```
x[y]
```

```
## [1] -0.1711196 -0.1499089 0.1173934
```

```
length(x)
```

```
## [1] 100
```

```
length(x[y])
```

```
## [1] 3
```

Exercise #1

Look at the output of the following code, what happens when you select with negative index?

```
x[3]
```

```
## [1] -0.1020015
```

```
x[-3]
```

```
## [1] -0.99307913 -0.17111961 -0.14990885 -0.21135557 -0.60078127 0.11739337
## [7] -0.21578390 -0.75237199 1.55294726 0.77671331 0.32967895 -1.89405513
## [13] -1.99349413 -0.22007207 0.19850126 0.56371430 -0.88079442 -0.64174738
## [19] -0.84737192 1.83233164 -0.53413321 0.30973709 0.01403343 -1.09327949
## [25] -0.07513613 1.52135816 0.53832708 -0.26987651 0.26677743 -0.70915890
## [31] -1.27608985 0.52580565 0.71607236 -0.60290030 0.14002073 -1.99917225
## [37] 0.42774676 -0.08996853 -0.23563275 -0.41061735 0.15146705 -2.26942347
## [43] -0.93500192 -1.41073822 -2.32958937 -0.15649706 -0.94767522 0.73027385
## [49] 0.41461206 0.65270923 0.64319947 -0.08089847 0.27330967 -0.08968567
## [55] -0.37876327 -1.15709419 0.58358139 -0.07484820 0.85993573 0.70481069
## [61] 0.66931003 0.38240296 0.01418130 1.05815293 -0.97596451 -0.01815480
## [67] -1.60150608 -0.29764517 -0.97331417 -1.53441965 2.09839585 0.95002228
## [73] 0.09173245 1.54627570 -1.42063598 0.73546857 0.12136776 1.03977273
## [79] -1.38440186 0.05776784 0.17496418 -0.19319262 0.06309704 0.52483624
## [85] 0.12168457 1.60986947 0.96609625 -0.89391965 -0.89837740 -1.56742899
## [91] 1.73821583 -1.21415287 1.44463281 -1.61656091 0.40198802 0.15775583
## [97] 0.57040588 0.07979506 0.43204664
```

ANSWERS: When a negative index is selected, the element at the selected index is removed from the vector.

Calculate summary statistics

We can calculate some summary statistics from the data. Here we set the option `collapse=T` for our code chunk, so that R Markdown will try to collapse all the source and output blocks from one code chunk into a single block.

```
mean(x)
## [1] -0.09494496
median(x)
## [1] -0.002060686

var(x) # sample variance
## [1] 0.8931129
sd(x) # sample standard deviation
## [1] 0.9450465

n <- length(x)
sum((x-mean(x))^2)/n # population variance
## [1] 0.8841817
mean(x^2)-(mean(x))^2 # population variance
## [1] 0.8841817
var(x)*(n-1)/n # population variance
## [1] 0.8841817
sd(x)^2*(n-1)/n # population variance
## [1] 0.8841817

summary(x)
##      Min.      1st Qu.      Median      Mean      3rd Qu.      Max.
## -2.329589 -0.776122 -0.002061 -0.094945  0.528936  2.098396
```

Exercise #2

Complete the following code to print out the maximum value and the minimum value in `x`.

ANSWERS

```
max(x)

## [1] 2.098396
```

```
min(x)

## [1] -2.329589
```

Exercise #3

The `sd()` function calculate the sample standard deviation.

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

How to calculate the population standard deviation?

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

ANSWERS

```
sd(x) # sample standard deviation s
## [1] 0.9450465

# calculate the population standard deviation
sd(x) * sqrt((n - 1) / n)
## [1] 0.9403094
```

Plot single data set

```
# generate random normals
set.seed(167) # feel free to change 167 to your lucky number
x <- rnorm(100) # use <- for assignment instead of =
```

Note that you can modify chunk options `fig.height` and `fig.width` to change the size of plots from a particular chunk.

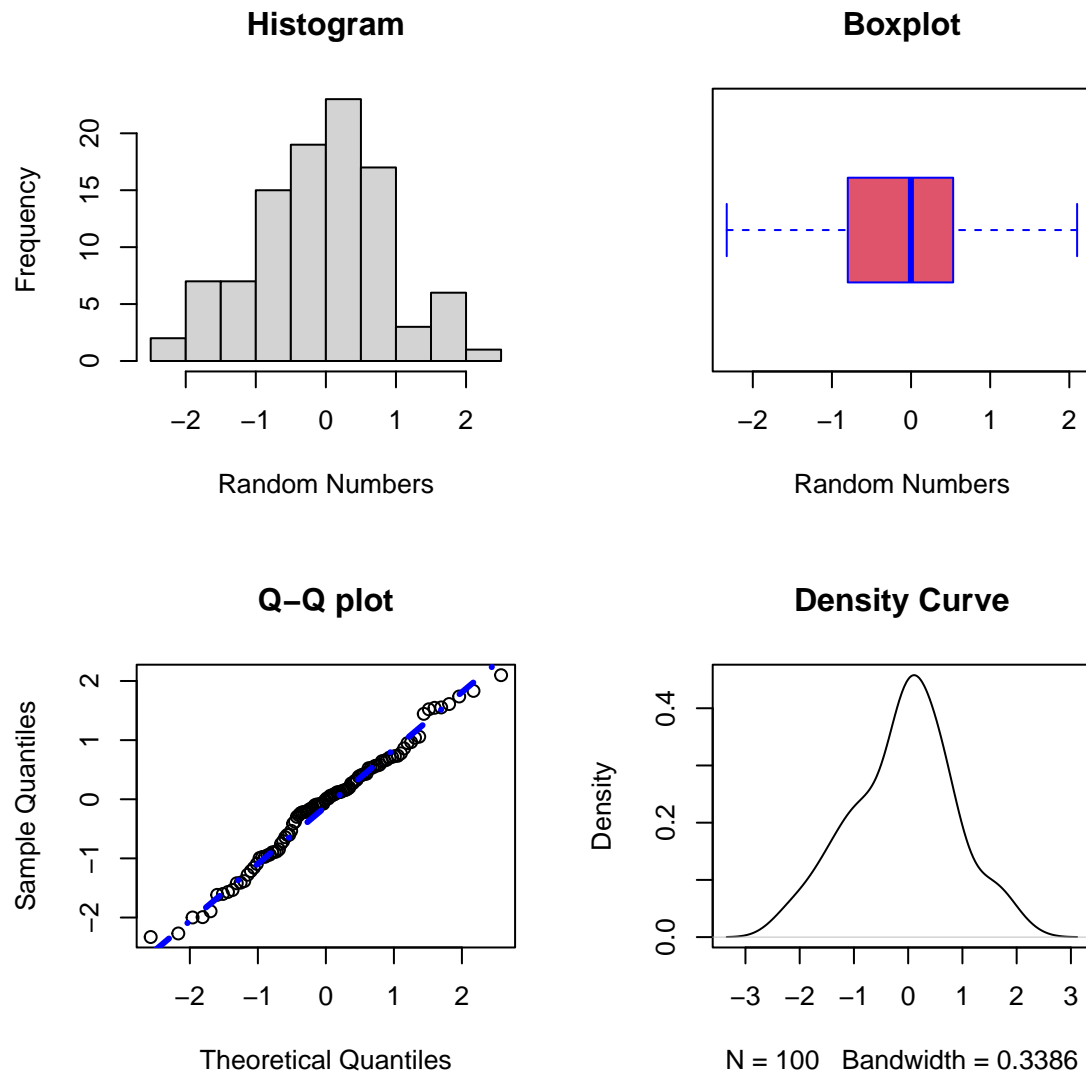
```
par(mfrow=c(2,2)) # 2-by-2 layout

# subfigure #1
hist(x, xlab="Random Numbers", main="Histogram")

# subfigure #2
boxplot(x, horizontal=T, col=2, border="blue", xlab="Random Numbers", main="Boxplot")

# subfigure #3
qqnorm(x, main="Q-Q plot")
qqline(x,lwd=3, lty=4, col="blue")

# subfigure #4
plot(density(x), main="Density Curve")
```

Plot a theoretical standard normal curve

```
z <- seq(-4, 4, length=100)
z
##      [1] -4.00000000 -3.91919192 -3.83838384 -3.75757576 -3.67676768 -3.59595960
##      [7] -3.51515152 -3.43434343 -3.35353535 -3.27272727 -3.19191919 -3.11111111
##     [13] -3.03030303 -2.94949495 -2.86868687 -2.78787879 -2.70707071 -2.62626263
##     [19] -2.54545455 -2.46464646 -2.38383838 -2.30303030 -2.22222222 -2.14141414
##     [25] -2.06060606 -1.97979798 -1.89898990 -1.81818182 -1.73737374 -1.65656566
##     [31] -1.57575758 -1.49494949 -1.41414141 -1.33333333 -1.25252525 -1.17171717
##     [37] -1.09090909 -1.01010101 -0.92929293 -0.84848485 -0.76767677 -0.68686869
##     [43] -0.60606061 -0.52525253 -0.44444444 -0.36363636 -0.28282828 -0.20202020
##     [49] -0.12121212 -0.04040404  0.04040404  0.12121212  0.20202020  0.28282828
##     [55]  0.36363636  0.44444444  0.52525253  0.60606061  0.68686869  0.76767677
##     [61]  0.84848485  0.92929293  1.01010101  1.09090909  1.17171717  1.25252525
```

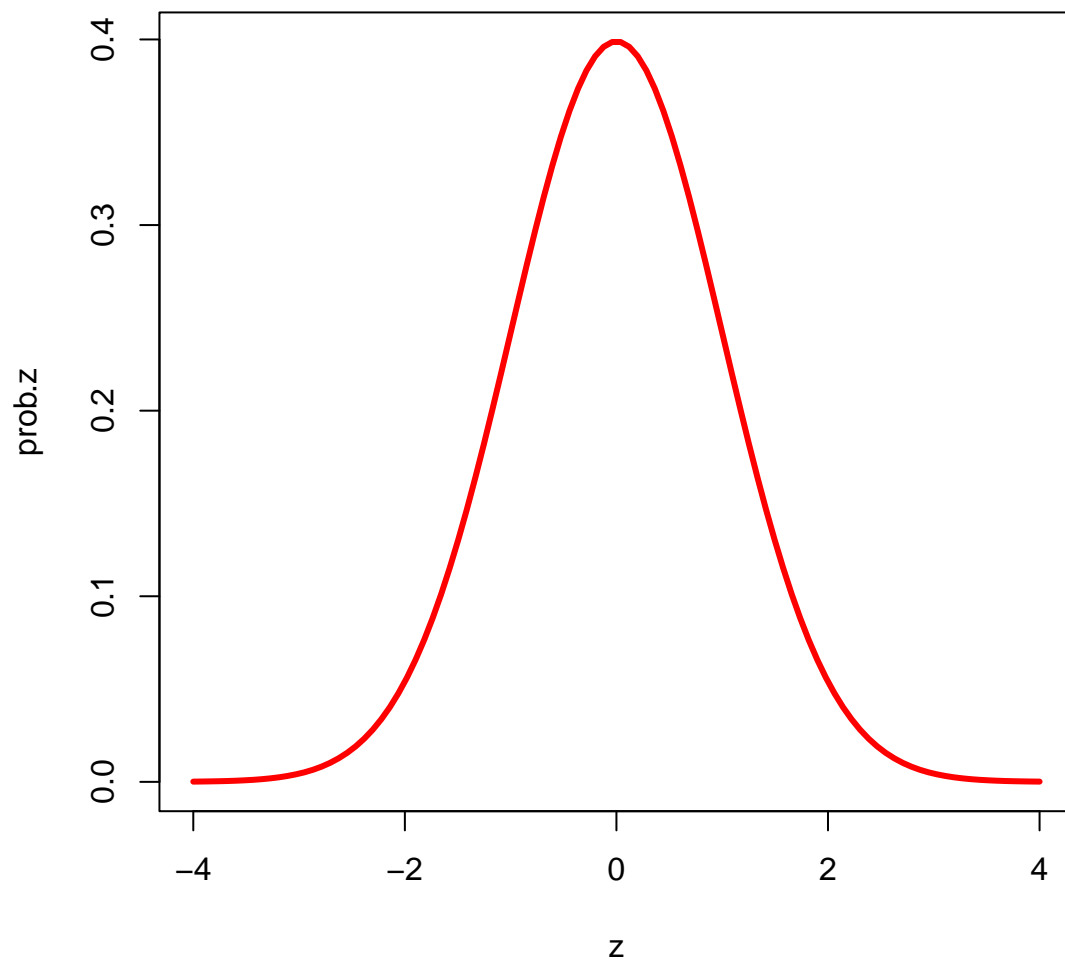
```
## [67] 1.33333333 1.41414141 1.49494949 1.57575758 1.65656566 1.73737374
## [73] 1.81818182 1.89898990 1.97979798 2.06060606 2.14141414 2.22222222
## [79] 2.30303030 2.38383838 2.46464646 2.54545455 2.62626263 2.70707071
## [85] 2.78787879 2.86868687 2.94949495 3.03030303 3.11111111 3.19191919
## [91] 3.27272727 3.35353535 3.43434343 3.51515152 3.59595960 3.67676768
## [97] 3.75757576 3.83838384 3.91919192 4.00000000
```

```
prob.z <- dnorm(z) # calculate the standard normal densities
```

```
prob.z
```

```
## [1] 0.0001338302 0.0001842953 0.0002521381 0.0003427099 0.0004627846
## [6] 0.0006208623 0.0008275148 0.0010957722 0.0014415473 0.0018840898
## [11] 0.0024464615 0.0031560163 0.0040448664 0.0051503080 0.0065151783
## [16] 0.0081881065 0.0102236211 0.0126820683 0.0156292995 0.0191360817
## [21] 0.0232771927 0.0281301641 0.0337736510 0.0402854146 0.0477399263
## [26] 0.0562056185 0.0657418315 0.0763955298 0.0881978860 0.1011608535
## [31] 0.1152738702 0.1305008512 0.1467776382 0.1640100747 0.1820728700
## [36] 0.2008093962 0.2200325354 0.2395266587 0.2590507715 0.2783428081
## [41] 0.2971250031 0.3151102096 0.3320089800 0.3475371752 0.3614238299
## [46] 0.3734189738 0.3833010942 0.3908839312 0.3960223134 0.3986167793
## [51] 0.3986167793 0.3960223134 0.3908839312 0.3833010942 0.3734189738
## [56] 0.3614238299 0.3475371752 0.3320089800 0.3151102096 0.2971250031
## [61] 0.2783428081 0.2590507715 0.2395266587 0.2200325354 0.2008093962
## [66] 0.1820728700 0.1640100747 0.1467776382 0.1305008512 0.1152738702
## [71] 0.1011608535 0.0881978860 0.0763955298 0.0657418315 0.0562056185
## [76] 0.0477399263 0.0402854146 0.0337736510 0.0281301641 0.0232771927
## [81] 0.0191360817 0.0156292995 0.0126820683 0.0102236211 0.0081881065
## [86] 0.0065151783 0.0051503080 0.0040448664 0.0031560163 0.0024464615
## [91] 0.0018840898 0.0014415473 0.0010957722 0.0008275148 0.0006208623
## [96] 0.0004627846 0.0003427099 0.0002521381 0.0001842953 0.0001338302
```

```
plot(z, prob.z, lwd=3, col="red", type="l") # draw the theoretical normal curve
```



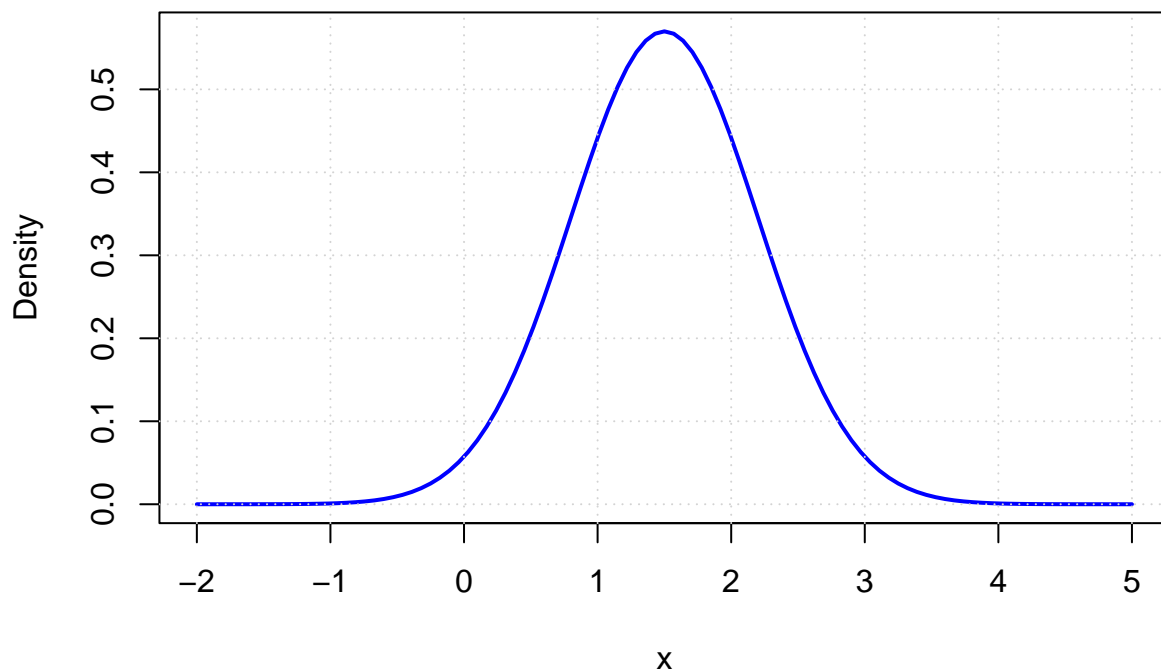
Exercise #4

Plot the probability density function of $N(\mu = 1.5, \sigma = 0.7)$.

ANSWERS

```
# I used curve to avoid hassle of seq(x)
curve(dnorm(x, mean = 1.5, sd = 0.7), from = -2, to = 5, lwd = 2, col = "blue",
      main = expression(paste("Probability Density Function for ", N(mu==1.5, sigma==0.7))),
      xlab = "x", ylab = "Density")
grid()
```

Probability Density Function for $N(\mu = 1.5, \sigma = 0.7)$



Can you draw these two normal curves, $N(\mu = 0, \sigma = 1)$ and $N(\mu = 1.5, \sigma = 0.7)$, in the same plot?

ANSWERS

```
curve(dnorm(x, mean = 0, sd = 1), from = -4, to = 4,
      col = "blue", lwd = 2,
      ylim = c(0, 0.75),
      xlab = "x", ylab = "Density",
      main = "Two N Curves")

curve(dnorm(x, mean = 1.5, sd = 0.7), add = TRUE,
      col = "red", lwd = 2)

grid()
```

Two N Curves

