

# 007 - Sampling Distributions

EPIB 607 - FALL 2020

Sahir Rai Bhatnagar  
Department of Epidemiology, Biostatistics, and Occupational Health  
McGill University

`sahir.bhatnagar@mcgill.ca`

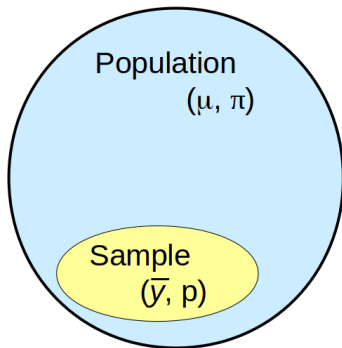
slides compiled on September 18, 2020





# Parameters and Statistics

- **Parameter:** An unknown numerical constant pertaining to a population/universe, or in a statistical model.
  - ▶  $\mu$ : population mean                       $\pi$ : population proportion
- **Statistic:** A numerical quantity calculated from a sample. The empirical counterpart of the parameter, used to *estimate* it.
  - ▶  $\bar{y}$ : sample mean                       $p$ : sample proportion



# Examples

## Proportions:

- Proportion of Earth's surface covered by water
- Proportion who saw a medical doctor last year
- Proportion of Québécois who don't have a family doctor

## Means:

- Mean depth in  $n$  randomly selected ocean locations
- Mean household size in  $n$  randomly selected households.
- Median number of persons under-5 in a sample of  $n$  households

# Samples must be random

- The validity of inference will depend on the way that the sample was collected. If a sample was collected badly, no amount of statistical sophistication can rescue the study.
- Samples should be **random**. That is, there should be no systematic set of characteristics that is related to the scientific question of interest that causes some people to be more likely to be sampled than others. The simplest type of randomization selects members from the population with equal probability (a uniform distribution).
- When conducting a study, it is always better to seek statistical advice sooner rather than later. Get a statistician involved at the *planning* stage of the study... by the analysis stage, it may be too late!

# Samples must be random - No cheating!

## Do not cheat by

- Taking 5 people from the same household to estimate
  - ▶ proportion of Québécois who don't have a family doctor
  - ▶ who saw a medical doctor last year
  - ▶ average rent
- Sampling the depth of the ocean only around Montreal to estimate
  - ▶ proportion of Earth's surface covered by water

# Collecting data takes effort

## In general

- The larger the sample  $\rightarrow$  the more accurate the estimate (if sampling is done correctly)

## CAVEAT

- Collecting more data takes effort and money!
- We will also soon discover the curse of the  $\sqrt{n}$

# Collecting data takes effort

## In general

- The larger the sample  $\rightarrow$  the more accurate the estimate (if sampling is done correctly)

## CAVEAT

- Collecting more data takes effort and money!
- We will also soon discover the curse of the  $\sqrt{n}$





# Sampling Distributions

- Given a sample of  $n$  observations from a population, we will be calculating estimates of the population mean, proportion, standard deviation, and various other population characteristics (parameters)
- Prior to obtaining data, there is uncertainty as to which of all possible samples will occur
- Because of this, estimates such as  $\bar{y}$  (the sample mean) will vary from one sample to another

# Sampling Distributions

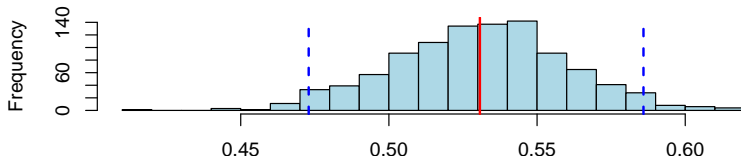
- The behavior of such estimates in many samples of equal size is described by what are called **sampling distributions**
- DVB definition: If we could see all the statistics (means, proportions, ect.) from all possible samples (Chapter 18, page 432)

# Sampling distribution of correlations<sup>1</sup>

Lets create a pseudo population from the 595 observations by sampling **with replacement**, and calculate the correlation. Lets repeat this process 1000 times:

```
library(oibiostat); data("famuss"); B <- 1000; N <- 595
R <- replicate(B, {
  dplyr::sample_n(famuss, size = N, replace = TRUE) %>%
    dplyr::summarize(r = cor(height, weight)) %>%
    dplyr::pull(r)
})
```

Distribution of samples of size 595



<sup>1</sup>from 004-exploring-data-2

# Why are sampling distributions important?

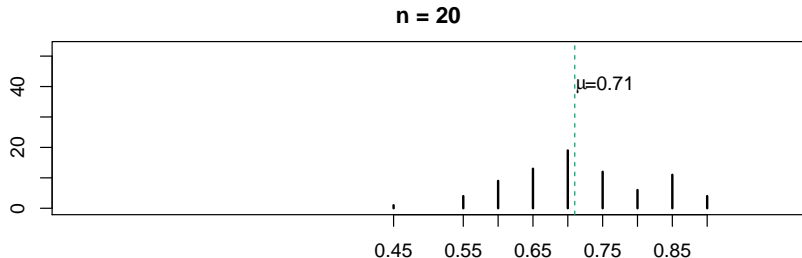
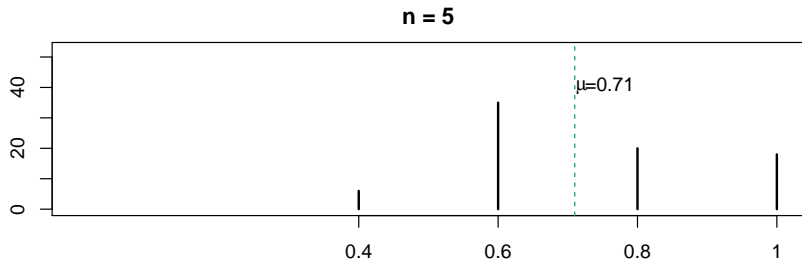
- Modeling how sample statistics vary from sample to sample is one of the most powerful ideas we'll see in this course.
- A sampling distribution *model* for how a sample statistics varies from sample to sample allows us to quantify that variation and to talk about how likely it is that we'd observe a sample statistic in any particular interval.
- Thus, they are used in confidence intervals for parameters. Specific sampling distributions (based on a null value for the parameter) are also used in statistical tests of hypotheses.

# Exercise 1: How Deep is the Ocean?

- We will get a sense of what a sampling distribution is in Exercise 1
- **CAVEAT:** This is a luxury using a toy example. In actual studies, we only get one shot!

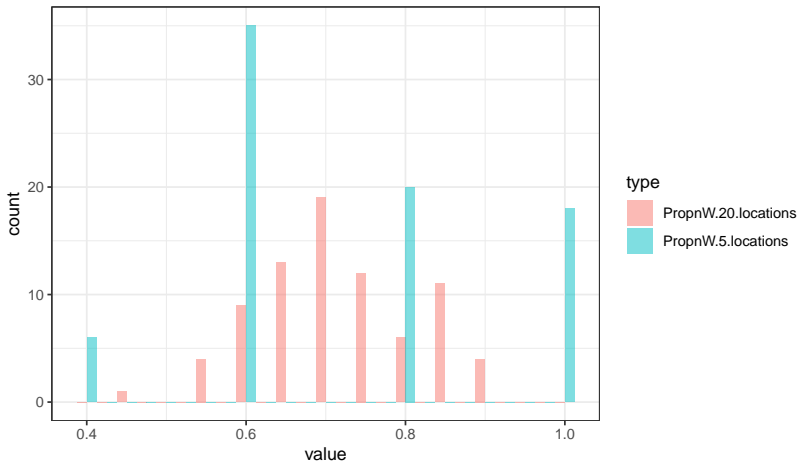


# Sampling distribution: proportion covered by water



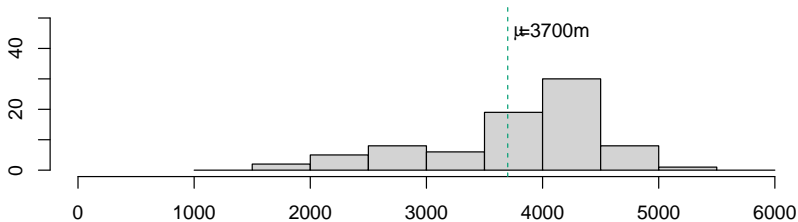


# Sampling distribution: proportion covered by water

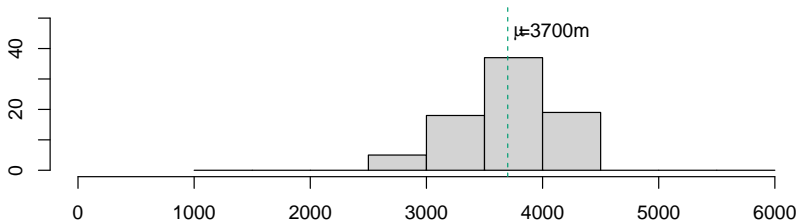


# Sampling distribution: mean depth of the ocean

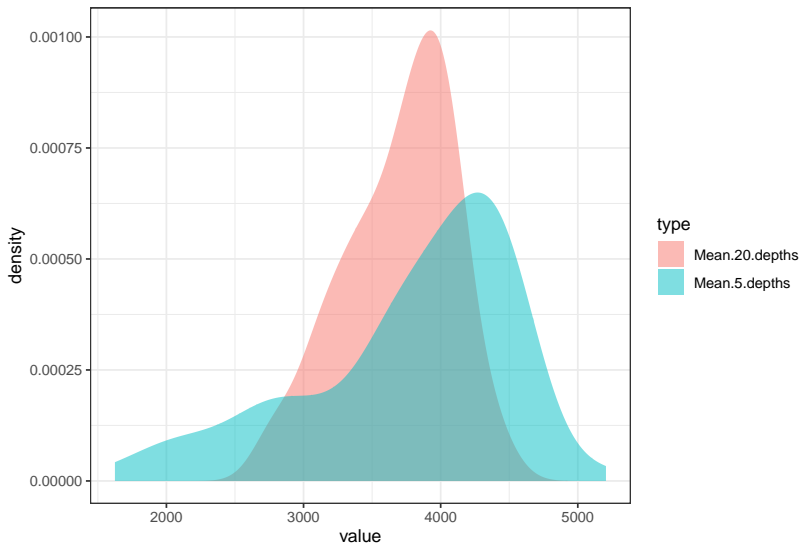
**n = 5**



**n = 20**



# Sampling distribution: mean depth of the ocean





# The Normal (Gaussian) distribution

What is it?

- A distribution that describes continuous (numerical) data
- Can also be used to approximate discrete data distributions
- Range is (technically) infinite, though the probability of seeing very large or very small values is extremely tiny
- Fully described by only two parameters, the mean and variance ( $\mu$  and  $\sigma^2$ )
- **NOTE:** R use the short-hand:  $X \sim \mathcal{N}(\mu, \sigma)$ , denoting the normal distribution as a function of the mean and *standard deviation*. This is not standard; many texts instead write  $X \sim \mathcal{N}(\mu, \sigma^2)$ . Be careful of this!

# The Normal (Gaussian) distribution

Carl Gauss was a German mathematician who developed a number of important advances in statistics such as the method of least squares.



**Figure:** The Deutsche Bundesbank issued Deutsche Mark banknotes in 15 different denominations, including this 10 Deutsche Marks banknote featuring Carl Friedrich Gauss.

# The Normal distribution

Where do Normal data come from?

- Natural processes
  - ▶ Blood pressure
  - ▶ Height
  - ▶ Weight
- “Man-made” (or derived)
  - ▶ Binomial (proportion) and Poisson (count) data are approximately Normal under certain conditions
  - ▶ Sums and means of random variables (Central Limit Theorem)
  - ▶ Data can sometimes be made to look Normal via transformations (squares, logs, etc)

# The Normal distribution

For Normal data, we can use the ~~Gaussian tables~~ **R** to answer the questions:

- What is the probability that a single observation  $X$  is
  - ▶ greater than  $X^*$ ?
  - ▶ less than  $X^*$ ?
  - ▶ between  $X_L^*$  and  $X_U^*$ ?
- That is, we can find out information about the percent distribution of  $X$  as a function of thresholds  $X^*$ , or  $X_L^*$  and  $X_U^*$ .
- We can also use the ~~Normal tables~~ **R** to find out information about thresholds  $X^*$  that will contain particular percentages of the data. I.e., we can find what threshold values will
  - ▶ Exclude the lower  $\omega^*$ % of a population
  - ▶ Exclude the upper  $\omega^*$ % of a population
  - ▶ Contain the middle  $\omega^*$ % of a population



# The Normal distribution

We can use ~~the Gaussian tables~~ R to answer these questions **no matter what the values of**  $\mu$  and  $\sigma^2$ .

That is, the % of the Normal distribution falling between  $X_L^* = \mu - m_1\sigma$  and  $X_U^* = \mu + m_2\sigma$  where  $m_1, m_2$  are any multiples **remains the same** for any  $\mu$  and  $\sigma$ .

How so??

Because we can **standardize** any  $X \sim \mathcal{N}(\mu, \sigma)$  to find  $Z \sim \mathcal{N}(0, 1)$

# The Normal distribution

An illustration using IQ scores, which we presume have a  $\mathcal{N}(100, 13)$  distribution of scores.

**Q1:** What percentage of scores are **above** 130?

Two steps:

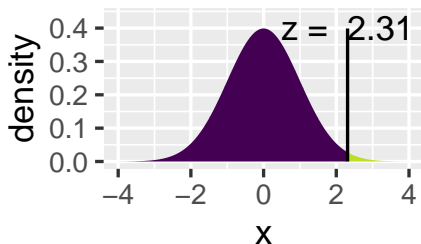
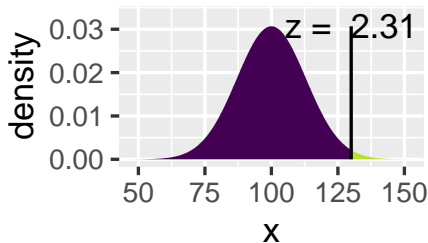
1. Change of location from  $\mu_X = 100$  to  $\mu_Z = 0$
2. Change of scale from  $\sigma_X = 13$  to  $\sigma_Z = 1$

Together, this gives us

$$Z = \frac{X - \mu_X}{\sigma_X} = \frac{130 - 100}{13} = 2.31$$

# The Normal distribution

The position of  $X=130$  in a  $\mathcal{N}(100, 13)$  distribution is the same as the place of  $Z = 2.31$  on the  $\mathcal{N}(0, 1)$ , which we call the **standardized** Normal distribution (or Z-distribution).



# The Normal distribution

How are the values in the Normal tables found?

Normal density:

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp \frac{-(x - \mu)^2}{2\sigma^2}$$

Probabilities found by integration (area under the Normal curve):

$$P(a \leq x \leq b) = \int_a^b \frac{1}{\sqrt{2\pi}\sigma} \exp \frac{-(x - \mu)^2}{2\sigma^2} dx$$

# The Normal distribution

(The percent above  $X = 130$ ) = (% above  $Z = 2.31$ ) = 1.04%

How do we know this? We look at the lower tail probability of 2.31 [i.e., the % below 2.31], and then subtract it from 1:

1.  $P(X < 130) = P(Z < 2.31) = 0.9896$
2.  $P(X > 130) = 1 - P(X < 130) = 0.0104$

So 130 is the 98.96<sup>th</sup> percentile of a  $\mathcal{N}(100,13)$  distribution.

# Reminder about percentiles and quantiles

- **Quantile**

- ▶ Any set of data, arranged in ascending or descending order, can be divided into various parts, also known as partitions or subsets, regulated by quantiles.
- ▶ Quantile is a generic term for those values that divide the set into partitions of size  $n$ , so that each part represents  $1/n$  of the set.
- ▶ Quantiles are not the partition itself. They are the numbers that define the partition.
- ▶ You can think of them as a sort of numeric boundary.

- **Percentile**

- ▶ Percentiles are quite similar to quantiles: they split your set, but only into two partitions.
- ▶ For a generic  $k$ th percentile, the lower partition contains  $k\%$  of the data, and the upper partition contains the rest of the data, which amounts to  $100 - k \%$ , because the total amount of data is  $100\%$ .
- ▶ Of course  $k$  can be any number between 0 and 100.

# More about percentiles and quantiles

- In class, we will find ourselves asking for the quantiles of a distribution.
- Percentiles go from 0 to 100
- Quantiles go from any number to any number
- Percentiles are examples of quantiles and you might find some people use them interchangeably (though this may not always be correct since quantiles can take on any value, positive or negative).
- **In particular**, R uses the term quantiles.
- **In the previous example**, we saw that  $P(Z < 2.31) = 0.9896$ . In R, 2.31 is called the quantile .

# The Normal distribution

(The percent above  $X = 130$ ) = (% above  $Z = 2.31$ ) = 1.04%

But wait!! The standard Normal is symmetric about 0, so we can do this another way... The % **above** 2.31 is equal to the % **below** -2.31:

$$\begin{aligned}P(X > 130) &= P(Z > 2.31) \\&\Rightarrow P(Z > 2.31) = P(Z < -2.31) \\&\Rightarrow P(X > 130) = P(Z < -2.31) = 0.0104\end{aligned}$$

So 130 is the 98.96<sup>th</sup> percentile of a  $\mathcal{N}(130, 13)$  distribution. What is the 1.04<sup>th</sup> percentile?

Transform from  $Z = -2.31$  back to  $X$ :

$$X = \sigma Z + \mu = 13(-2.31) + 100 = 69.97.$$

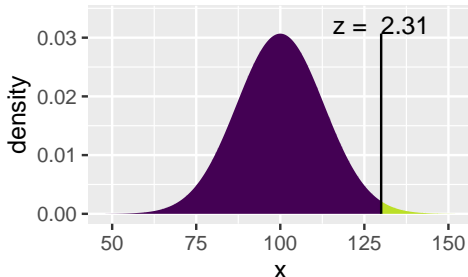


# For probabilities we use *pnorm*

```
stats::pnorm(q = 130, mean = 100, sd = 13)
```

```
## [1] 0.99
```

```
mosaic::xpnorm(q = 130, mean = 100, sd = 13)
```



```
## [1] 0.99
```

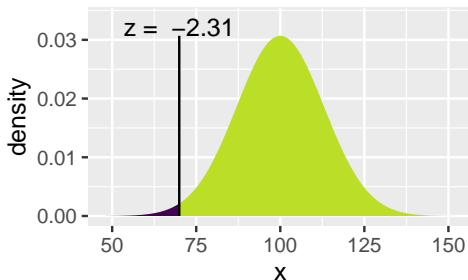
- `pnorm` returns the integral from  $-\infty$  to  $q$  for a  $\mathcal{N}(\mu, \sigma)$
- `pnorm` goes from *quantiles* (think *Z* scores) to probabilities

# For quantiles we use *qnorm*

```
stats::qnorm(p = 0.0104, mean = 100, sd = 13)
```

```
## [1] 70
```

```
mosaic::xqnorm(p = 0.0104, mean = 100, sd = 13)
```



```
## [1] 70
```

- *qnorm* answers the question: What is the Z-score of the  $p$ th percentile of the normal distribution?
- *qnorm* goes from *probabilities* to quantiles

# The Normal distribution

Q2: What is the probability of seeing an IQ score **as extreme as** (think highly unusual) 130?

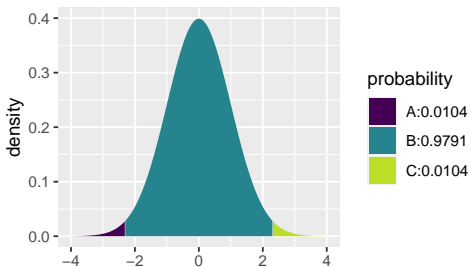
1. Again, we find that  $X = 130$  is the same percentile of the IQ Normal distribution as  $Z = 2.31$  is of the standard Normal.
2. To see what scores are as extreme, we want to know the probability that  $Z > 2.31$  or that  $Z < -2.31$ .
3. As we saw previously,  $P(Z > 2.31) = P(Z < -2.31) = 0.0104$ , so the probability of seeing an IQ as extreme or more so than 130 is  $2 \times 0.0104 = 0.0208$ .

# Finding tail probabilities

```
# lower.tail = TRUE is the default
stats::pnorm(q = -2.31, mean = 0, sd = 1, lower.tail = TRUE) +
stats::pnorm(q = 2.31, mean = 0, sd = 1, lower.tail = FALSE)

## [1] 0.021
```

```
mosaic::xpnorm(q = c(-2.31, 2.31), mean = 0, sd = 1)
```



```
## [1] 0.01 0.99
```

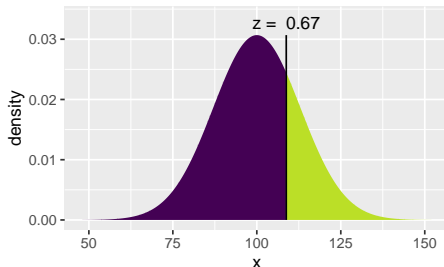
# The Normal distribution

Q3: What is the 75<sup>th</sup> percentile of the IQ scores distribution?

We now have to reverse the sequence of steps:

- **Ask yourself:** What Z value corresponds to a probability of 0.75? Should you use `pnorm` or `qnorm`?

```
mosaic::xqnorm(p = 0.75, mean = 100, sd = 13)
```



```
## [1] 109
```

This tells us that 75% of the IQ scores fall below 108.8.

# Empirical Rule or 68-95-99.7% Rule

In any normal distribution with mean  $\mu$  and standard deviation  $\sigma^2$ :

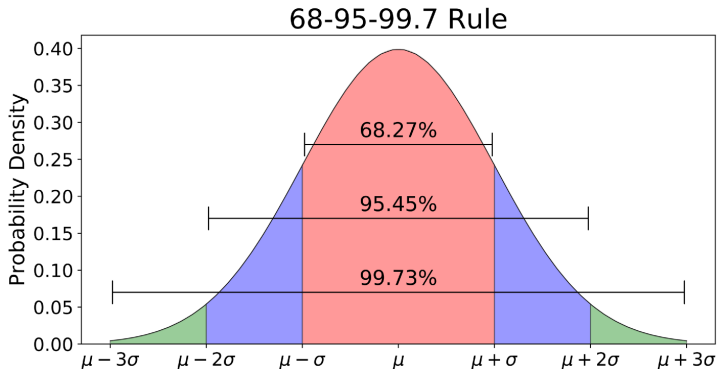
- Approximately 68% of the data fall within one standard deviation of the mean.
- Approximately 95% of the data fall within two standard deviations of the mean.
- Approximately 99.7% of the data fall within three standard deviations of the mean.

# Demo of Empirical Rule

```
pacman::p_load(mosaic)
pacman::p_load(manipulate)

mNorm <- function(mean = 0, sd = 1) {
  lo <- mean - 5 * sd
  hi <- mean + 5 * sd
  manipulate(
    xpnorm(c(A,B), mean, sd, verbose = FALSE, invisible = TRUE),
    A = slider(lo, hi, initial = mean - sd),
    B = slider(lo, hi, initial = mean + sd)
  )
}
mNorm(mean = 0, sd = 1)
```

# Empirical Rule or 68-95-99.7% Rule





# Properties of Normal random variables

Special properties of the Normal distribution:

- If  $Y$  is a Normal random variable, then so is  $a + bY$ .
- If  $X$  and  $Y$  are two Normal random variables, then  $X + Y$  is a Normal random variable. What is the mean and variance of this new random variable?
- If  $X$  and  $Y$  are two Normal random variables and  $\rho_{XY} = 0$  (correlation between  $X$  and  $Y$ ), then  $X$  and  $Y$  are independent.

# Properties of Normal random variables

Let  $Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma)$ , and let each  $Y_i$  be independent of the others.  
(think simple random sample)

Then  $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$  has what distribution?

- The sum of Normal random variables is Normal, so  $\bar{Y}$  is a Normal random variable.
- $E(\bar{Y}) = \frac{1}{n} \sum_{i=1}^n E(Y_i) = \frac{1}{n} \sum_{i=1}^n \mu = \mu.$
- $Var(\bar{Y}) = Var(\frac{1}{n} \sum_{i=1}^n Y_i) = \frac{1}{n^2} \sum_{i=1}^n Var(Y_i) = \sigma^2/n.$
- Standard Error of  $\bar{Y} = \sqrt{Var(\bar{Y})} = \sigma/\sqrt{n}$

# Session Info

```
R version 4.0.2 (2020-06-22)
Platform: x86_64-pc-linux-gnu (64-bit)
Running under: Pop!_OS 19.10

Matrix products: default
BLAS:   /usr/lib/x86_64-linux-gnu/openblas-pthread/libblas.so.3
LAPACK: /usr/lib/x86_64-linux-gnu/openblas-pthread/liblapack.so.3

attached base packages:
[1] tools      stats      graphics  grDevices  utils      datasets  methods
[8] base

other attached packages:
[1] mosaic_1.7.0      Matrix_1.2-18     mosaicData_0.20.1 ggformula_0.9.4
[5] ggstance_0.3.4    lattice_0.20-41   oibiostat_0.2.0   NCStats_0.4.7
[9] FSA_0.8.30        forcats_0.5.0     stringr_1.4.0     dplyr_1.0.2
[13] purrr_0.3.4       readr_1.3.1       tidyr_1.1.2       tibble_3.0.3
[17] ggplot2_3.3.2     tidyverse_1.3.0   knitr_1.29

loaded via a namespace (and not attached):
[1] fs_1.5.0           lubridate_1.7.9    httr_1.4.2         backports_1.1.9
[5] R6_2.4.1           DBI_1.1.0          colorspace_1.4-1   withr_2.2.0
[9] tidyrselect_1.1.0  gridExtra_2.3      leaflet_2.0.3      curl_4.3
[13] compiler_4.0.2     cli_2.0.2          rvest_0.3.6        xml2_1.3.2
[17] ggdendro_0.1.22    labeling_0.3       mosaicCore_0.8.0   scales_1.1.1
[21] digest_0.6.25      foreign_0.8-79     rio_0.5.16         pkgconfig_2.0.3
[25] htmltools_0.5.0    dbplyr_1.4.4       highr_0.8          htmlwidgets_1.5.1
[29] rlang_0.4.7        readxl_1.3.1       rstudioapi_0.11    farver_2.0.3
[33] generics_0.0.2     jsonlite_1.7.1     crosstalk_1.1.0.1  zip_2.1.1
[37] car_3.0-9          magrittr_1.5       Rcpp_1.0.5         munsell_0.5.0
[41] fansi_0.4.1        abind_1.4-5        lifecycle_0.2.0    stringi_1.5.3
[45] carData_3.0-4      MASS_7.3-53        plyr_1.8.6         grid_4.0.2
[49] blob_1.2.1         ggrepel_0.8.2      crayon_1.3.4       haven_2.3.1
[53] splines_4.0.2      hms_0.5.3          pillar_1.4.6       reprex_0.3.0
[57] glue_1.4.2         evaluate_0.14      data.table_1.13.0  modelr_0.1.8
[61] vctrs_0.3.4        tweenr_1.0.1       cellranger_1.1.0   gtable_0.3.0
[65] polyclip_1.10-0    assertthat_0.2.1   TeachingDemos_2.12 xfun_0.17
[69] ggforce_0.3.2      openxlsx_4.1.5     broom_0.7.0        viridisLite_0.3.0
[73] ellipsis_0.3.1
```