. Within R a list is a structure that can combine objects of different types.

- A Boolean indicator of whether or not they have insurance. • Name • Name • Amount of bill due

An Example Database

We first consider a patient database where we want to store their

A ni sted gnitteð - f ysd Intro to R Programming for Biostatistics

- A list is a type of vector called a recursive vector. - Other vectors are atomic vectors . A list is actually a vector but it does differ in comparison to the other types of vectors which we have been using in this class.

Creating Lists

Introduction to R

· We will learn how to create and work with lists in this section.

Lists

Lists

Creating the Same List

"sizil"=abom)" - yec'so" - 'jis.s "signA" -> [["amen"]]jis.s ?7 -> [["bawo"]]jis.s TRUE -> [["abomenuzul"]]jis.s

[7] LENE
[7] _VVB6F9,,
[7] _VVB6F9,,
2096G

Types of Information

We then have 3 types of information here:

To create a list of one patient we say

Single Patient

lesinemun ·

· logical.

- · Our database would then be a list of all of these individual lists.

Lists of Lists

. We could then create a list like this for all of our patients.

[7] IMRE ## 27U2NL9UCG ## [7] "22" ## ## [7] "Vu8659" ## [7] "Vu8659" a <- list(name="Angela", owed="75", insurance=TRUE) 6

List Operations

- . With vectors, arrays and matrices, there was really only one way to index them.
- · However with lists there are multiple ways:

. There is another easy way to create this same list $% \left\{ 1,2,\ldots,n\right\}$. This also allows us to help with indexing as we will see below.

Notice that unlike a typical vector this prints out in multiple parts.

Buixəpul

Double vs Single Brackets

"character" [1] ##
cpez([[]])
"siagna" [1] ##
[[τ]]e

List Indexing

name"]]	,]]e	
"sisgna" [í	##	
π	[]]e	
"efagov" [1	##	
au	eu\$e	
"s19gnA" [I	##	

Double vs Single Brackets

- . With the single bracket we are returned another list.
- . With the double bracket we are returned an element in the original class of what kind of data we entered.

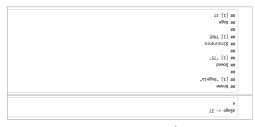
Double vs Single Brackets

- All of the previous are ways to index data in a list.
 Notice that in two of the above we used double brackets.
- Next we see the difference between double and single brackets.

- \cdot Depending on your goals you may want to use single or double brackets.

Adding and Subtracting Elements

With a list we can always add more information to it.



Double vs Single Brackets

## [1] "list"	
cjss2(9[j])	
## \$nane ## [1] "Angela"	
[1]	

Adding and Subtracting Elements

In order to delete an element from a list we set it to NULL.

	6
amen\$	пп
"siegnA" [i]	
	пп
\$Ţuznusuce	пп
aumī [1]	пп
	пп
age\$	пп
ΔZ [τ]	пп

List Components and Values

In order to know what kind of information is included in a list we can look at the $\ensuremath{\textit{nones0}}$ function

[1] "character"

Unlisting

"age" "aonenuent" "amsn" [1] ## (e)saueu

Applying Functions to Lists

- Just like arrays and matrices we can use an apply() function.
 Specifically we have lapply() and sapply() functions for lists.
- With the original ${\it appiy0}$ function we could specify whether the function was applied to either the rows or the columns.

· If your list contained all numerical elements than the class would be numerical. · If There is Character data in the original list that unlisted everything will be in character format.

BnitsilnU

- \cdot With the case of lists both functions are applied to elements of the list.

Applying Functions to Lists

```
[2]] ***

| TE 96 SE 76 SE 77 SE 78 SE 61 SE 71 SE 71 SE 71 SE 72 SE 76 SE 78 
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       #Wumber list
n <- list(l:5, 6:37)
n
```

name insurance

To find the values of things we could go ahead and unlist them

Recursive Lists

· This is because we can actually have lists within lists. · Before it was mentioned that a list is a recursive vector.

- Applying Functions to Lists
- . With this list we see that we have two separate vectors of numbers included.
- . Then let us see the results of either using lapply() and sapply()

Recursive Lists

For example let us go back to our patient data.

```
patients <- list(a,s)
```

([t]] ## ([c]] ## ([c]]

```
S'TZ 0'E [T] ##
```

Applying Functions to Lists

Apply Functions an Lists

Recursive Lists

```
96 [1]
10 [3]
10 [4]
11 [4]
12 [4]
12 [4]
13 [4]
14 [5]
15 [6]
16 [5]
16 [6]
16 [6]
17 [6]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18 [7]
18
```

. While the sapply() function returns a vector of the medians.

- The **lapply()** function returns a list with the median of each of the original lists.

Example with Output of a List



Final Motes on Lists

- . It is important to remember how we can call these features of lists.
- You almost never want to use the generated output from R. . Many of you will want to use R for model building and regressions.
- . For example $\ensuremath{\text{R}}$ does not automatically return the confidence intervals with a regression.

Example with Output of a List

- So R just gave me the coefficients back but no other information.
 This means my knowledge of accessing lists is key.

"feubizen.ib" "febom"	"ap" "smn91"		"səufev.bəttt" [2]	
"rank"	"effects"	"residuals"	"stnsisitheos" [1]	##

Final Motes on Lists

- · This is why we take the time to discuss how to search what is in a list and how to access it.

- . The output from most regression functions in R is actually a list.
- What this means is I can extract the elements from the list that I want in order to build tables
 that display the exact information that I want it to.

Example with Output of a List

- $\,$ I can see that R actually has a lot more information that they did not display for me.
- Next I consider a function where it summarizes the information from this model

Example with Output of a List

x <- rnorm(500,10, 3) y <- 3*x + rnorm(500, 0, 2)

Conclusion of Lists

Example with Output of a List

- $\,$ R has so much information about regression that is never even displayed unless I dig deeper.
- Understanding lists and accessing information means you can output custom tables that look much more professional than what R gives you.

Example with Output of a List

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27% 2.4%)

(2.27%)

(2.27% 2.4%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

(2.27%)

DataFrames in R

Dataframe

- . With statistics we are most likely to use the data structure called a data frame.
- This is similar to a matrix in appearance however we can have multiple types of data in it like \dots
- Each column must contain the same type of data or R will most likely default to character for that column.
- It is very important that you become proficient in working with data frames in order to fully understand data analysis.

Example with Output of a List

names ages from mext.appt
TRUE 2016-9-23
A Short of the control of the cont

patients <- chind(patients, next.appt)
patients

("erbrod2" ,"SelgghA")s -> samen (26,5%)s -> samen (26,5%)s -> sage (28,5%)s -> samen (1,28%)s -> samensat. (1,28%)s -> samensat. (28,5%)s -> samensat. (2

Adding Rows or Columns

. We may wish to add rows or columns to our data.

()puiqo -

1 Angela 252 trrurence
2 Anondra 36 TRUE
3 Shondra 74 ANS

l <- c(names="Liu lie", age=45, insurance=TRUE) rbind(patients, l)

Adding Rows or Columns

Wenching in `[<-.factor`(`*tmp*', ri, value = "Liu Jie"): invalid factor ## Zevel, MB generated

batieut we coniq just do the following For example we are to add another

- We can do this with:

Next appointments $\label{eq:condition} \text{mext.appt} \leftarrow c("09)23/2016", "04/14/2016", "02/25/2016")$

Finally if we decided to then place another column of data in we could

Adding Rows or Columns

patlents\$names <- as.character(patlents\$names)
patlents <- rbind(patlents, 1)
patlents **Adding Rows or Columns**

We usually create a data frame with vectors. **Creating Data Frames**

Bnixabnl

- · We can also access data by using both column and row information • head() function allows us to access the first rows of the data frame.
 - · Indexing is the same as that for matrices.

Accessing Data Frames

In order to best consider accessing of data frames we will use some built in data from $\ensuremath{\ensuremath{R_{\text{\tiny L}}}}$

"x92" "Sex" "Sex"

Variables Included in Titanic

library(datasets) titanic -- data.frame(Titanic)

Bnixabnl

** [1] Onlid Onlid Child # accessing age information titanic[,3]

[1] Onld Onld Child Onld Onld Child Child Child Onld Child Onld Child Onld Child Onld Child C

BrimeM bring Brixabril

- This means we can access data with a column or row number
- · More importantly we can use the name.
- . For large data frames accessing by a name is key.

head(titanic) ## Class Sex Age Survived Freq ## 1 1st Male Child No 0 ## 2 Zhd Male Child 0 0 titanic[l:2,] Preview Into Data

Class Sox Age Survived Freq ### The Fred Child No 8 ### Crew Holle Child No 8 ### Crew Holle Child No 9 ### Crew Holle Child No 9 ### Crew Holle Child No 9

Replacing Values

. Perhaps we were not pleased the decimal places and want to have this as a percentage.

· We can overwrite the values and change this.

				(=(=	TURN TO	١nı	
d ⁻ Auns	pəu∃	panţnung	93A	xəs	CJass		##
000000.0	0	ON	ситла	Male	1st	τ	##
000000.0	0	ON	ситла	Male	puz	ζ	##
981065°T	32	ON	сиття	Male	bns	٤	##
000000.0	0	ON	сиття	Male	CLGM	Þ	##
	000000.0 000000.0	000000,0 0 000000,0 0 081002,1 28	881862,1 2E ON	Child No 0 0.000000 Child No 0 0.0000000	Sex Age Survived Freq survive Male Child No 9 0.808080 Male Child No 3 1.590186	Class Sex Age Survived Freq surv. Draft Set Fit and Fit of 6 6,00000000 and Male Child No 35 1,500186	1 1st Male Child No 0 0.000000 2 2nd Male Child No 0 0.00000000000000000000000000000000

Further Indexing

We could ask for information by using the factors that we have as well

Z6T	SŁ	ÞΤ	۷S	0	ΣŢ	ττ	S	029	48 E	⊅ST	811	0	32	0	0	[τ]	##	
					["bə.	ч.,	"ale",	d==:	kə\$\$	oţue	171]oţue	2772			male.	
								ØÞ1	I ZS	τ	S	Þ	81	τø	0	[1] ##	
			["bəu	ч. '	,,151	==	:sseTD\$	ţoţu	2171]ɔṛu	671:					isa£‡ isa£‡	

Our New Variables

. We can ask R to calculate this and add it to our data.

Seldsing Wew Variables

- Suppose we not only want to know the frequency of survival but the proportion

TE2T [T] ##
(male.freq)
SZE [I] ##
(pani.zasi.trati)mus

A ni zəlddiT

zelddi**T**

Previously we have worked with data in the form of

- sasia · · Vectors
- Arrays
- · Dataframes
- $\label{eq:constraint} \mbox{titanic$Freq/sum(titanic$Freq)} \mbox{ read(titanic$freq)}$

Non-Standard Names

· "Tibbles" are a new modern data frame. **z**elddiT

- · It removes many of the outdated features. · It keeps many important features of the original data frame.

names(data.frame('crazy name' = 1)) ## [1] "crazy.name" names(tibble('crazy name' = 1)) ## [2] "crazy name"

- · This works similar to as.data.frame().

Coercing into Tibbles

- A tibble can be made by coercing as_tibble().

It works efficiently.

· A tibble can have non-standard variable names. - No more worry of characters being automatically turned into strings.

- To use this refer to these in a backtick.

· A tibble never changes the input type.

Compared to Data Frames

- It only recycles vectors of length 1.

Coercing into Tibbles

```
## puris: microseconds

### participation of mathem of mathematical mat
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      microbenchmark::microbenchmark(
as_tibble(l),
as.data.frame(l)
)
                                                                                                                                                                                                                                                                                                                                                                      l <- replicate(26, sample(100), simplify = FALSE)
names(l) <- letters
```

```
$3.5.5 undersolv (Lidywerse)

1. Tyer Citable(s = 1.5.) y = lite(1.5.) (1.10)

1. Tyer Citable(s = 1.5.) y = lite(1.5.) (1.10)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.10, 1.20))

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.10, 1.20))

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.10, 1.20))

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.10, 1.20))

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.10, 1.20))

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Citable(s = 1.5.) y = lite(1.5., 1.20)

1. Tyer Cit
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 column-Lists
```

Snittesduc

Zibbles vs Data Frames

There are a couple key differences between tibbles and data frames.

· Subsetting.

[["x"]]tb -

[[t]]≯p -

. We can also use a pipe which we will learn about later.

 $\boldsymbol{\cdot}$. We can index a tibble in the manners we are used to

[["x"]]. %<% hb -

Subsetting

```
XTD

ABACH22.0 A22ETI2.0 CELIZ28.0 ELIGAGE, D ESGYIZA.0 [1] ##

ABACH22.0 A22ETI2.0 CELIZ28.0 ELIGAGE, D ESGYIZA.0 [1] ##

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]**

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[[1]]*

[
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             df <- tibble(
x = runif(5),
y = rnorm(5)
)</pre>
```

```
tibble(
a = lubridate::now() + runif(las) * 86489,
b = lubridate::toda() + runif(las) * 30,
c = llas,
d = runif(las),
)
)
)
```

Tibbles only print the first 10 rows and all the columns that fit on a screen. - Each column displays Its data type.

Subsetting

```
As and abschez, a perettra, a celtrean, a estancer, a estancer, a estancer, a estancer, a estancer, and a abschez, a perettra, a celtrean, a estancer, a estancer, a estancer, a estancer, a estancer, and a perettra a celtra and a celtra an
```

Printing

You will not accidentally print too much.

Printing

```
| 0564/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0564/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-00-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U199''' # ##

| 0664/1E1'0 01 | 07-000 066 U
```