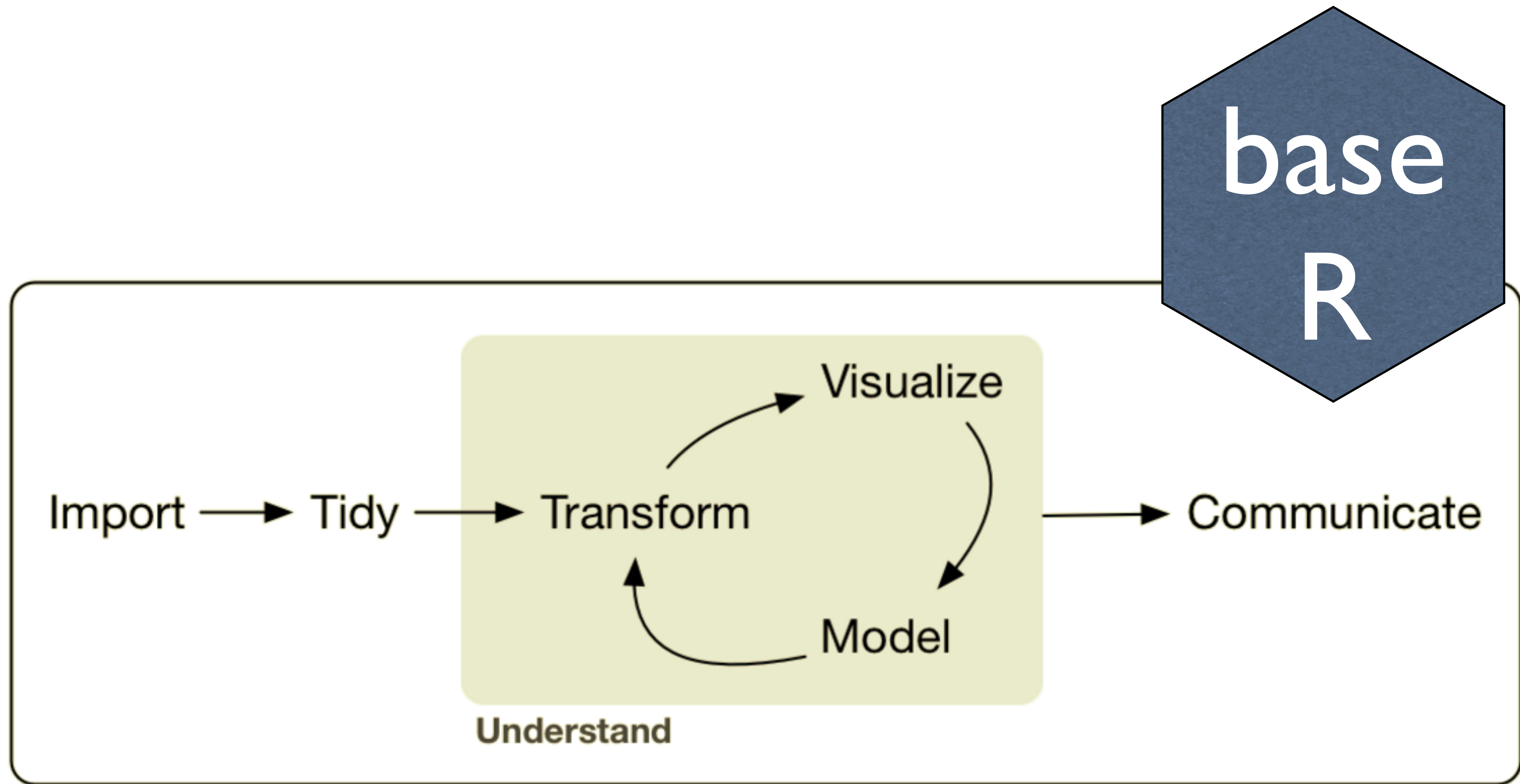


# More on Data Types and Syntax — transitioning to programming in R

# Data Types in R

1. Vectors
2. Matrices
3. Data types
4. Coercion



**Program**

From *R for Data Science* by Hadley Wickham and Garrett Grolemund.

# Your Turn 1

Look at the R object

`WorldPhones` (by typing its name in your notebook or the Console and hitting enter).

What is inside of `WorldPhones`?

A digital timer with a black border and white background, displaying the time 00:30 in a black, segmented font.

# WorldPhones

	N.Amer	Europe	Asia	S.Amer	Oceania	Africa	Mid.Amer
1951	45939	21574	2876	1815	1646	89	555
1956	60423	29990	4708	2568	2366	1411	733
1957	64721	32510	5230	2695	2526	1546	773
1958	68484	35218	6662	2845	2691	1663	836
1959	71799	37598	6856	3000	2868	1769	911
1960	76036	40341	8220	3145	3054	1905	1008
1961	79831	43173	9053	3338	3224	2005	1076

You can save more than a single number in an object by creating a *vector*, *matrix*, or *array*.

# Vectors

# Your turn

How many dimensions does a vector have?

1 2 3 4 5 6

00:30

How many dimensions does a vector have?





# vectors

Combine multiple elements into a one dimensional array.

Create with the `c` function (for "concatenate").

```
vec <- c(1, 2, 3, 10, 100)
```

```
vec
```

# Your turn

What happens in your Environment  
when you run this code?

In your Notebook?

```
vec <- c(1, 2, 3, 10, 100)
```

```
vec
```

A digital timer with a black border and white background, displaying the time 00:30 in a black, segmented font.

# vectors

Combine multiple elements into a one dimensional array.

Create with the `c` function (for "concatenate").

```
vec <- c(1, 2, 3, 10, 100)
```

```
vec
```

```
# 1  2  3 10 100
```



# Matrices

Uncommon, but good to know

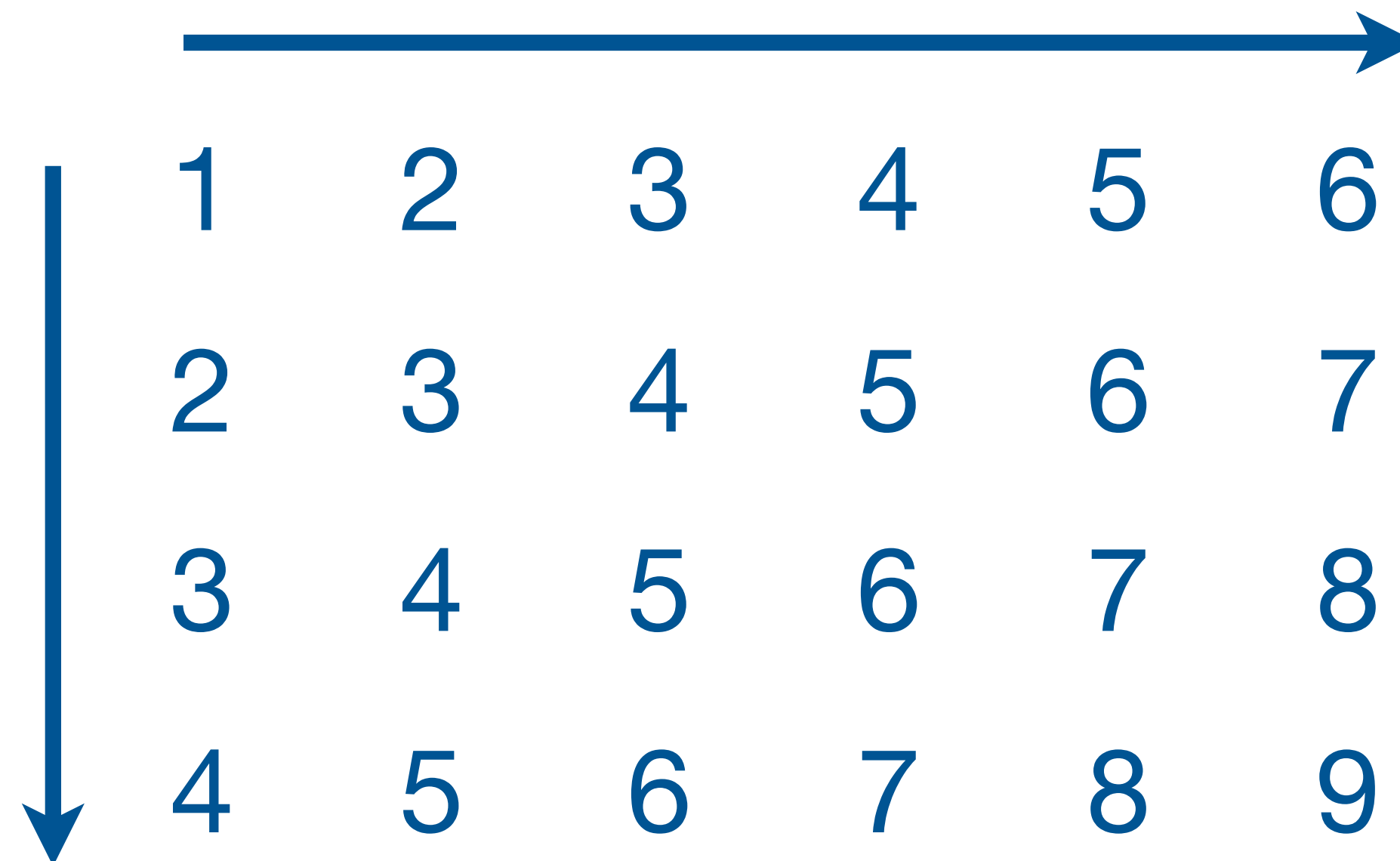
# Your turn

How many dimensions does a matrix have?

1	2	3	4	5	6
2	3	4	5	6	7
3	4	5	6	7	8
4	5	6	7	8	9

00:30

How many dimensions does a matrix have?



# Your turn

The matrix below is named **M**.  
What is the value of **M<sub>34</sub>**?

0	1	2	3	4	5
6	7	8	9	10	11
12	13	14	15	16	17
18	19	20	21	22	23



00:30

The matrix below is named **M**.  
What is the value of **M**<sub>34</sub>?



0	1	2	3	4	5
6	7	8	9	10	11
12	13	14	15	16	17
18	19	20	21	22	23




The matrix below is named **M**.  
What is the value of  **$M_{34}$** ?



	0	1	2	3	4	5
	6	7	8	9	10	11
	12	13	14	15	16	17
	18	19	20	21	22	23

The matrix below is named **M**.  
What is the value of  **$M_{34}$** ?



0	1	2	3	4	5
6	7	8	9	10	11
12	13	14	15	16	17
18	19	20	21	22	23

# matrices

multiple elements stored in a two dimensional array.

Create with the `matrix` function.

```
mat <- matrix(c(1, 2, 3, 4, 5, 6), nrow = 2)
```

```
mat
```

```
#      [,1] [,2] [,3]
```

```
# [1,]  1  3  5
```

```
# [2,]  2  4  6
```

# matrices

Combine multiple elements into a two dimensional array.

Create with the `matrix` function.

```
mat <- matrix(c(1, 2, 3, 4, 5, 6), nrow = 2)
```

```
mat
```

```
#      [,1] [,2] [,3]
```

```
# [1,]  1   3   5
```

```
# [2,]  2   4   6
```



vector of elements to  
go in the matrix

```
matrix(c(1, 2, 3, 4, 5, 6), nrow = 2)
```

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

number of rows for  
matrix

```
matrix(c(1, 2, 3, 4, 5, 6), nrow = 2)
```

```
#      [,1] [,2] [,3]
```

```
# [1,]  1   3   5
```

```
# [2,]  2   4   6
```

```
matrix(c(1, 2, 3, 4, 5, 6), nrow = 3)
```

```
#      [,1] [,2]
```

```
# [1,]  1  4
```

```
# [2,]  2  5
```

```
# [3,]  3  6
```

R as a  
calculator  
(again)



# Math: element-wise

vec + 4

# 5 6 7 14 104

vec \* 4

# 4 8 12 40 400

vec \* vec

# 1 4 9 100 10000

The image shows a screenshot of the RStudio interface. The main editor window displays R code with several lines highlighted in green, indicating comments. The code includes a matrix creation function, a section header for math with vectors and matrices, and various operations on vectors and matrices. The environment pane on the right shows the current environment (Global Environment) with a list of objects: 'bechdel' (1794 obs. of 15 ...), 'mat' (num [1:2, 1:3] 1...), and 'vec' (num [1:5] 1 2 3 10...). The console at the bottom is empty.

```
27 matrix(c(1, 2, 3, 4, 5, 6), nrow = 3)
28 ```
29
30 ## Math with vectors and matrices
31
32 ```{r}
33 vec + 4
34 vec * 4
35 vec * vec # element-wise multiplication
36
37 vec %% vec # matrix multiplication (inner)
38 vec %%% vec # matrix multiplication (outer)
39
40 mat
41 t(mat) # transpose
42 ```
43
44 ## Arrays
45
46 ```{r}
47 array(c(1, 2, 3, 4, 5, 6), dim = c(2, 2, 3))
48 ```
```

Environment: Global Environment

Data	
bechdel	1794 obs. of 15 ...
mat	num [1:2, 1:3] 1...

Values	
vec	num [1:5] 1 2 3 10...

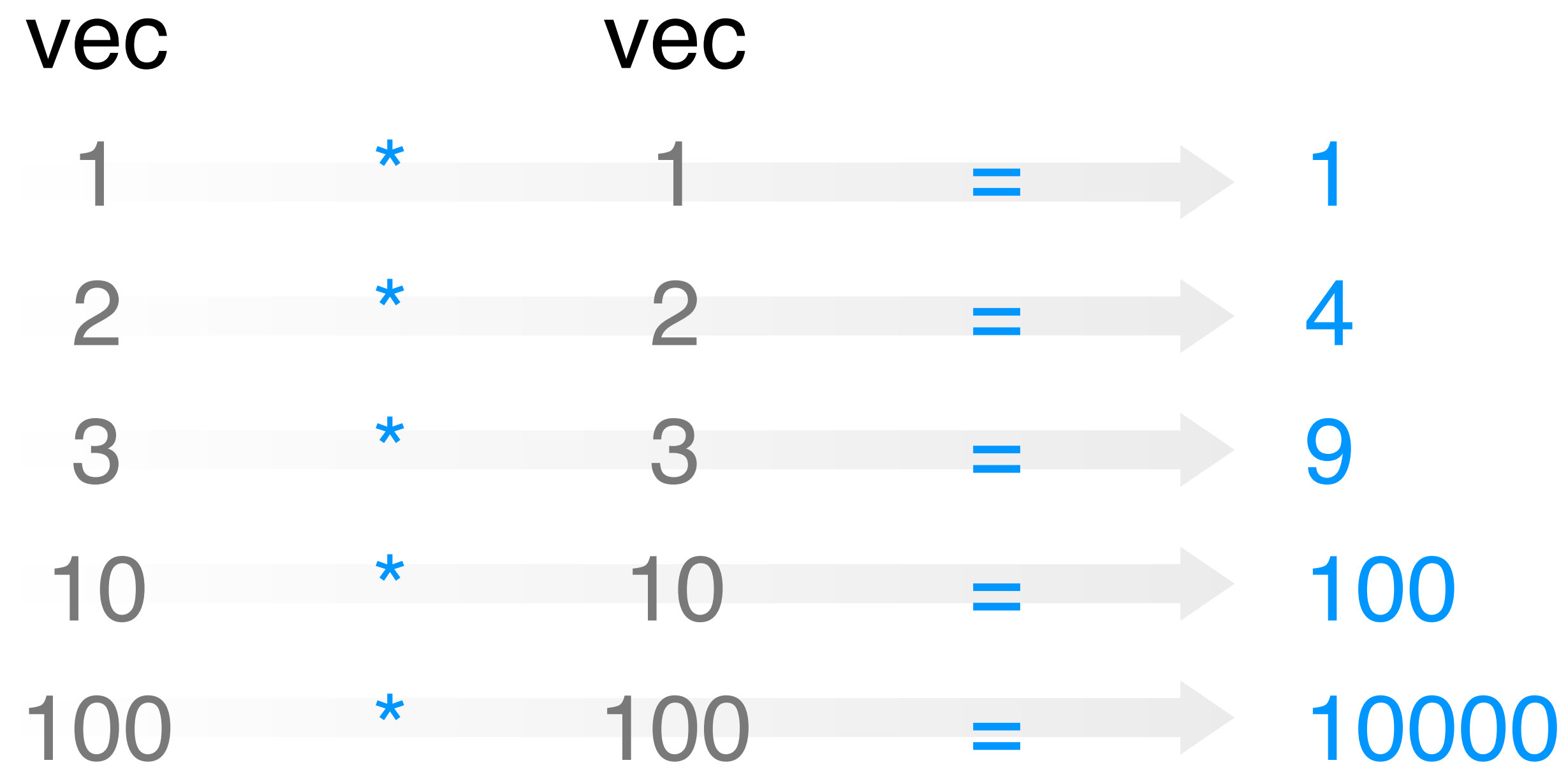
33:8 | Chunk 4 | R Markdown

Console

Green text indicates a code "comment," another way to document what you're doing. Comments aren't executed by R when you run a line.

vec \* vec

# 1 4 9 100 10000



# Matrix multiplication

```
vec %*% vec # inner
```

```
# [1]
```

```
# [1,] 10114
```

```
vec %o% vec # outer
```

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

```
# [1,] 1 2 3 10 100
```

```
# [2,] 2 4 6 20 200
```

```
# [3,] 3 6 9 30 300
```

```
# [4,] 10 20 30 100 1000
```

```
# [5,] 100 200 300 1000 10000
```

mat

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

t(mat)

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

# arrays

Combine multiple elements into an array that has three or more dimensions.

Create with the `array` function.

```
array(c(1, 2, 3, 4, 5, 6), dim = c(2, 2, 3))
```

# arrays

Combine multiple elements into an array that has three or more dimensions.

Create with the `array` function.

```
array(c(1, 2, 3, 4, 5, 6), dim = c(2, 2, 3))
```

Another uncommon  
structure

# Data types



# Warm up

	A	B	C	D
1	date	president	democrat	<u>unemploy</u>
2	Mar 31, 1968	Lyndon Johnson	TRUE	2709
3	Apr 30, 1968	Lyndon Johnson	TRUE	2740
4	May 31, 1968	Lyndon Johnson	TRUE	2938
5	Jun 30, 1968	Lyndon Johnson	TRUE	2883
6	Jul 31, 1968	Lyndon Johnson	TRUE	2768
7	Aug 31, 1968	Lyndon Johnson	TRUE	2686
8	Sep 30, 1968	Lyndon Johnson	TRUE	2689
9	Oct 31, 1968	Lyndon Johnson	TRUE	2715
10	Nov 30, 1968	Lyndon Johnson	TRUE	2685
11	Dec 31, 1968	Lyndon Johnson	TRUE	2718
12	Jan 31, 1969	Richard Nixon	FALSE	2692
13	Feb 28, 1969	Richard Nixon	FALSE	2712
14	Mar 31, 1969	Richard Nixon	FALSE	2758
15	Apr 30, 1969	Richard Nixon	FALSE	2713
16	May 31, 1969	Richard Nixon	FALSE	2816
17	Jun 30, 1969	Richard Nixon	FALSE	2868
18	Jul 31, 1969	Richard Nixon	FALSE	2816
19	Aug 31, 1969	Richard Nixon	FALSE	2816
20	Sep 30, 1969	Richard Nixon	FALSE	2816
21	Oct 31, 1969	Richard Nixon	FALSE	2816
22	Nov 30, 1969	Richard Nixon	FALSE	2816

What types of data appear in this spreadsheet?

# data types

Like Excel, Numbers, etc., R can recognize different types of data.

We'll look at four basic types:

- numbers
- character strings (text)
- logical
- factor

# numeric

Any number, no quotes.

Appropriate for math.

1 + 1

3000000

class(0.00001)

# "numeric"

# character

Any symbols surrounded by quotes.

Appropriate for words, variable names, messages, any text.

```
"hello"
```

```
class("hello")
```

```
# "character"
```

```
"hello" + "world"
```

```
# Error
```

```
nchar("hello")
```

```
# 5
```

```
paste("hello", "world")
```

```
# "hello world"
```

# Your turn

Which of these are numbers? **What are the others?** How can you tell?

1

"1"

"one"

00:30

# logical

TRUE or FALSE

R's form of binary data. Useful for logical tests.

```
3 < 4
```

```
# TRUE
```

```
class(TRUE)
```

```
# "logical"
```

```
class(T)
```

```
# "logical"
```

# factor

R's form of categorical data. Saved as an integer with a set of labels (e.g. levels).

```
fac <- factor(c("a", "b", "c"))
```

```
fac
```

```
# a b c
```

```
# Levels: a b c
```

```
class(fac)
```

```
# factor
```

Use great caution with  
factors



# Quiz

```
x <- c(1, 2, 3)
```

What is the difference between these?

x

"x"

00:30

<b>Type</b>	<b>Examples</b>
numeric	0, 1, -2, 3.1415, 0.0005
character	"Amelia", "Agree", "31"
logical	TRUE, FALSE
factor	a c c b Levels: a b c

# Your turn 2

Make a vector that contains the number 1, the letter R, and the logical TRUE.

What class of data is the vector?

```
vec <- c(1, "R", TRUE)
```

```
class(vec)
```

```
# "character"
```

```
vec
```

```
# "1" "R" "TRUE"
```

```
# What is R doing?
```

# Your turn

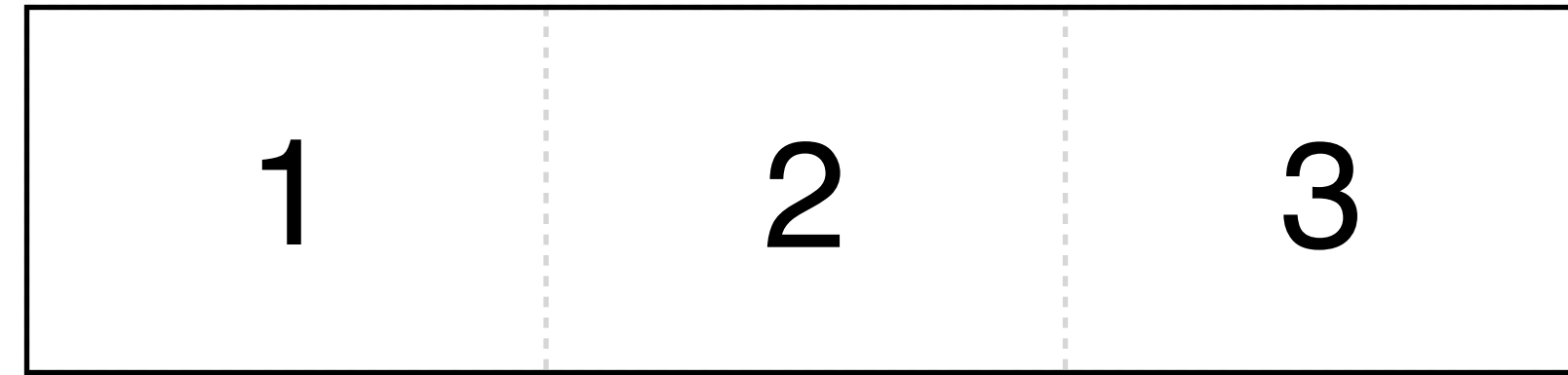
Another way to see the class of an object is in the Environment pane. Does the Environment agree with what you found using `class()`?

A digital timer with a black border and white background, displaying the time 00:30 in a black, segmented font.

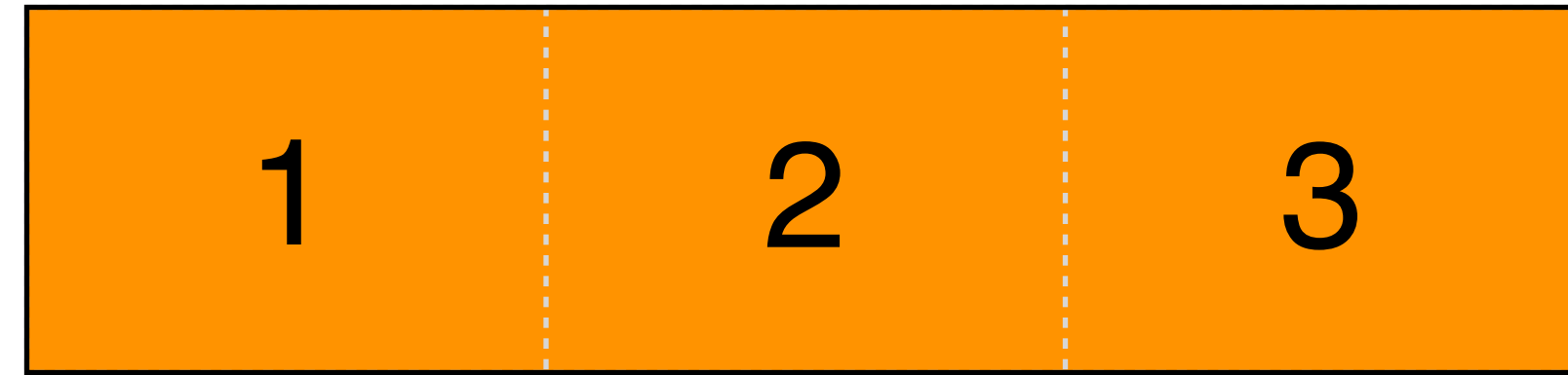
Vector



Vector



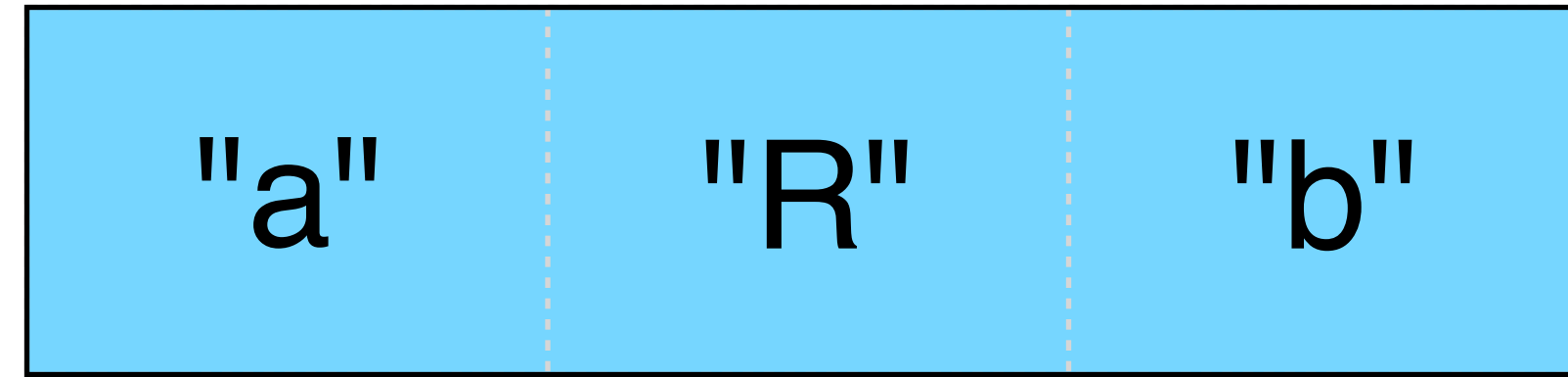
Vector



numeric

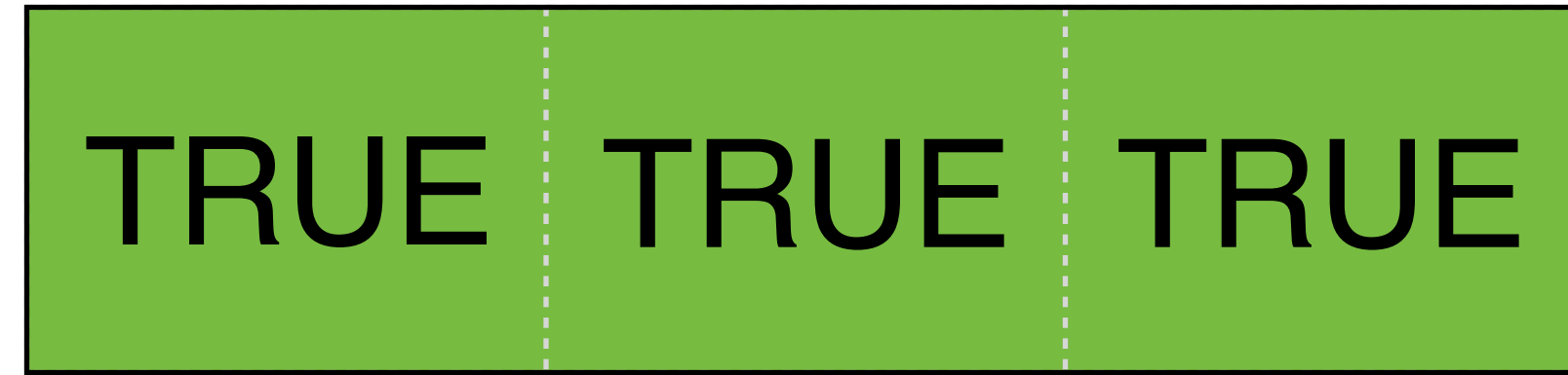


Vector



character

Vector



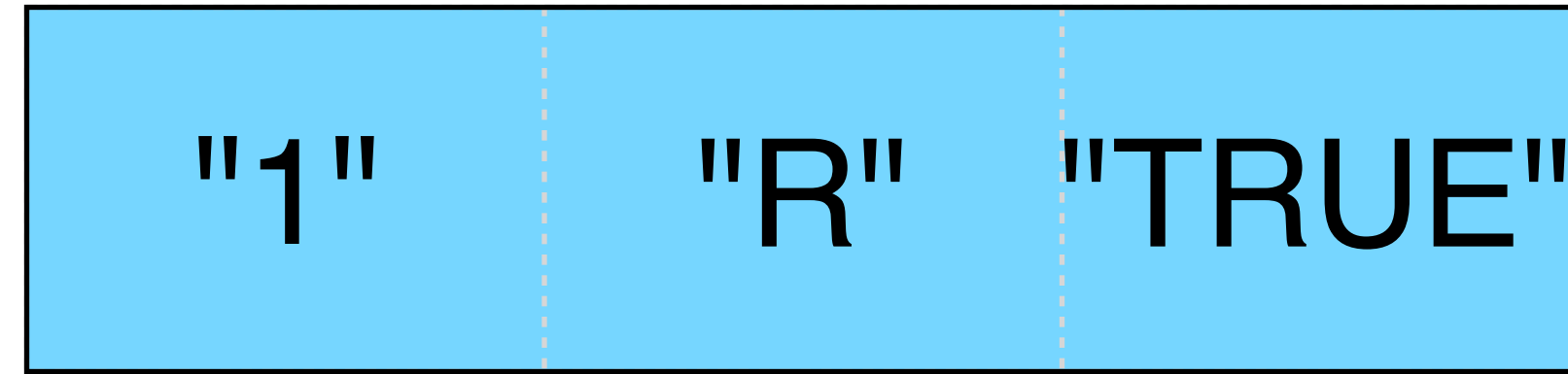
logical

Vector

1	"R"	TRUE
---	-----	------

?

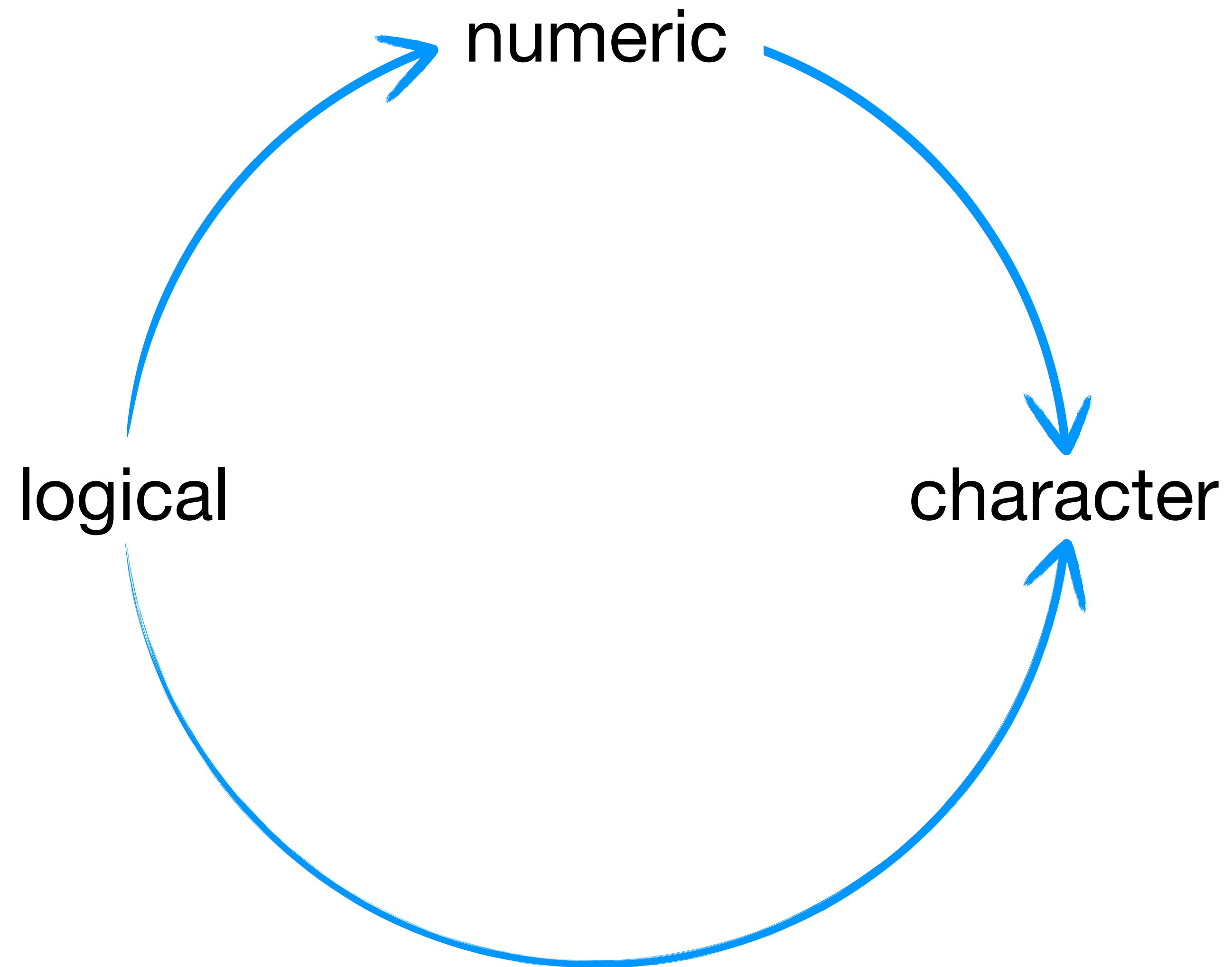
Vector



character

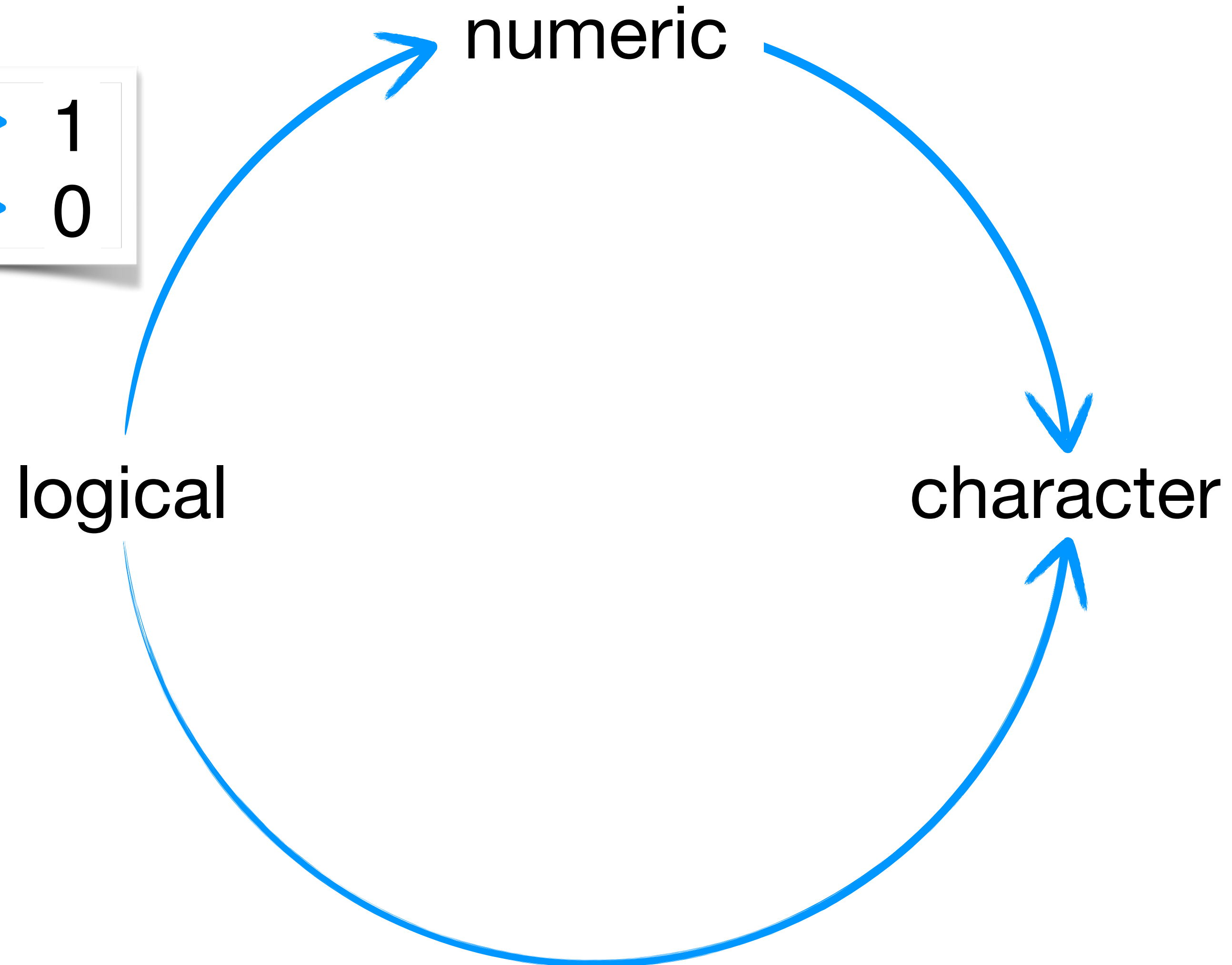
Coercion

# coercion

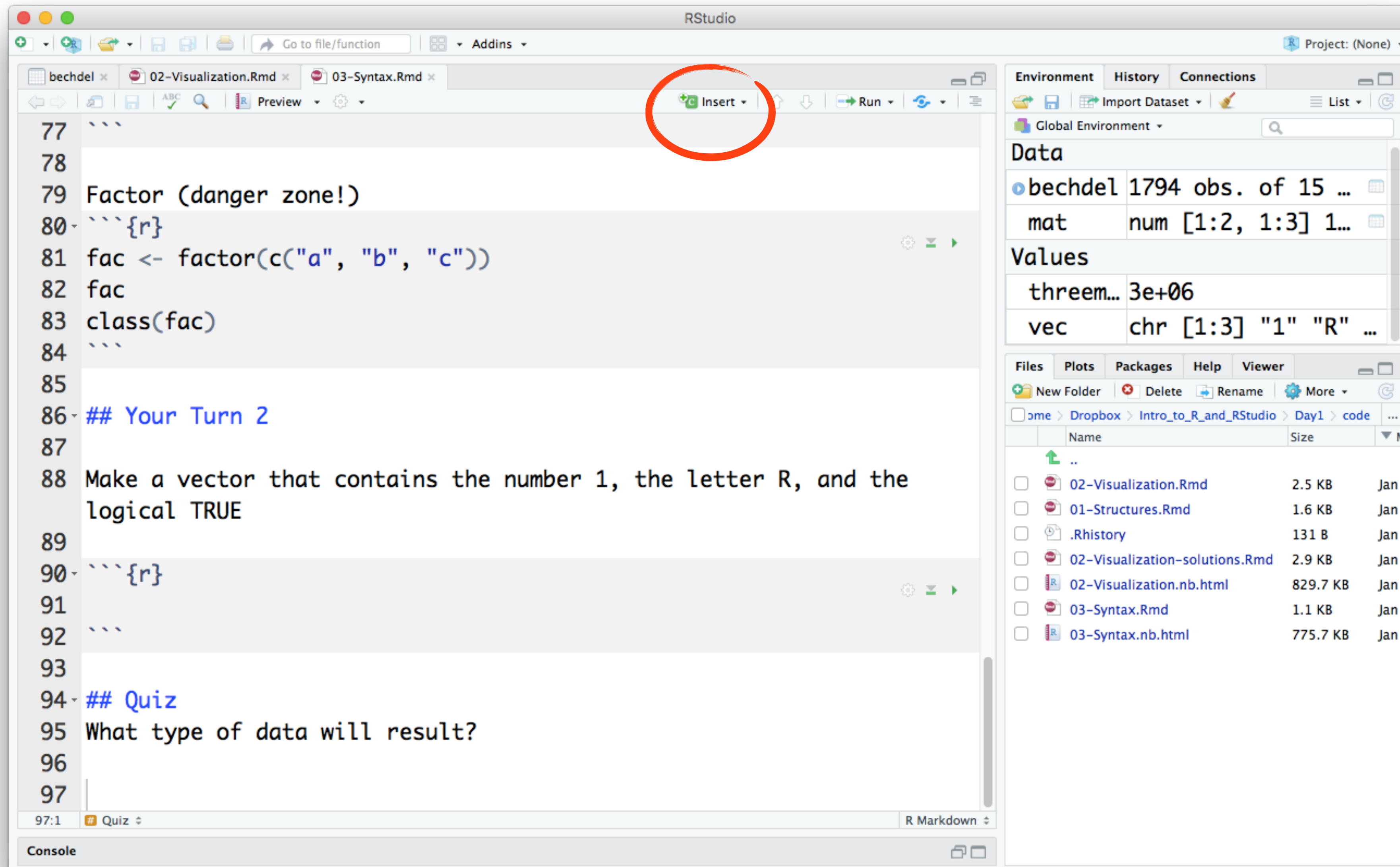


# coercion

TRUE → 1  
FALSE → 0



I'm going to give you a "quiz", and you might want to create your own chunk to try out some code. Use the Insert button to insert one



The screenshot shows the RStudio interface with the following content:

```
77 ```  
78  
79 Factor (danger zone!)  
80 ```{r}  
81 fac <- factor(c("a", "b", "c"))  
82 fac  
83 class(fac)  
84 ```  
85  
86 ## Your Turn 2  
87  
88 Make a vector that contains the number 1, the letter R, and the  
89 logical TRUE  
90 ```{r}  
91  
92 ```  
93  
94 ## Quiz  
95 What type of data will result?  
96  
97
```

The 'Insert' button in the toolbar is circled in red. The right-hand side of the interface shows the Environment pane with the following data:

Data	
bechdel	1794 obs. of 15 ...
mat	num [1:2, 1:3] 1...
Values	
threem...	3e+06
vec	chr [1:3] "1" "R" ...

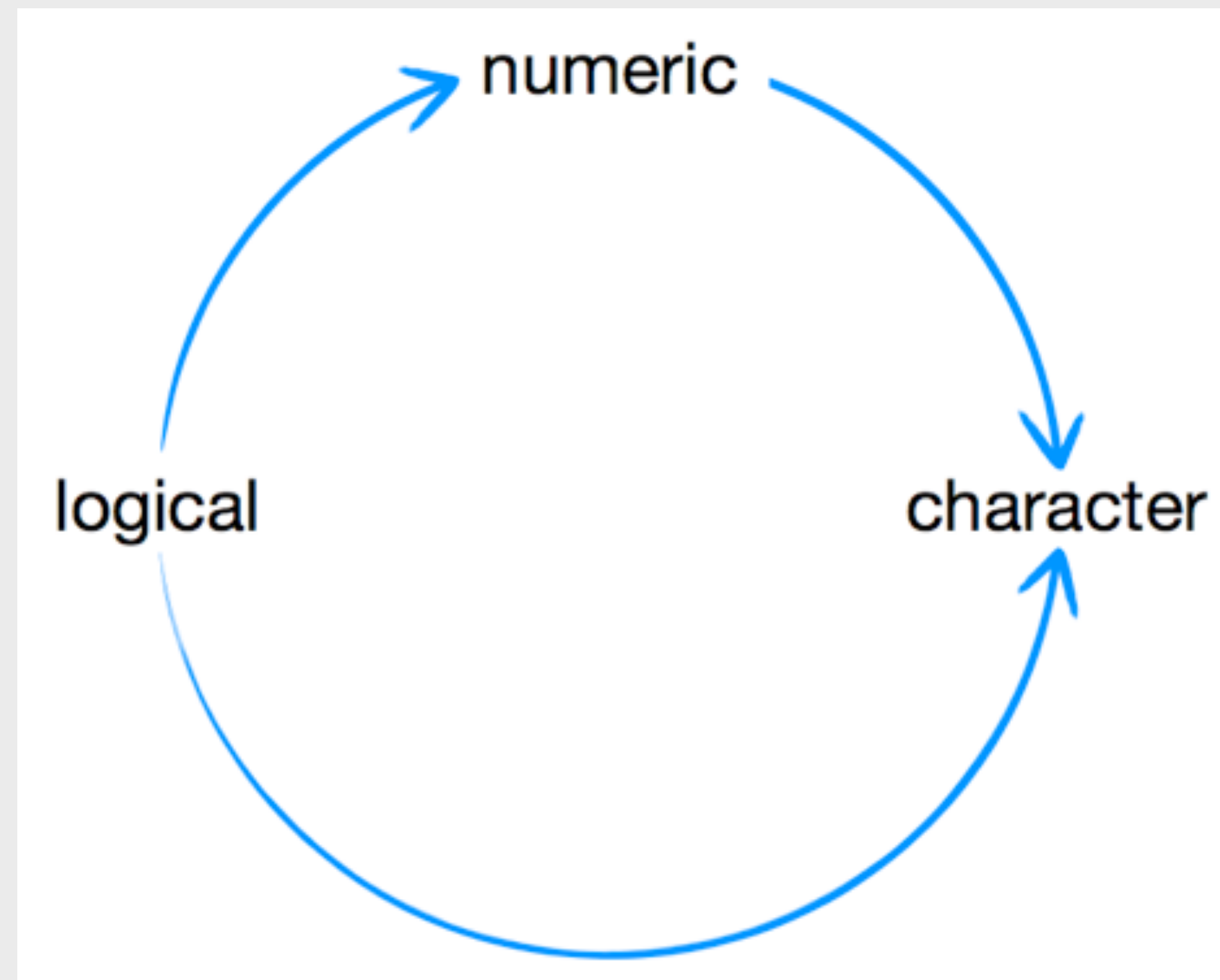
The Files pane shows the following files:

Name	Size	Modified
..		
02-Visualization.Rmd	2.5 KB	Jan
01-Structures.Rmd	1.6 KB	Jan
.Rhistory	131 B	Jan
02-Visualization-solutions.Rmd	2.9 KB	Jan
02-Visualization.nb.html	829.7 KB	Jan
03-Syntax.Rmd	1.1 KB	Jan
03-Syntax.nb.html	775.7 KB	Jan



# Quiz

What type of data will result?



`c(5, "two")`

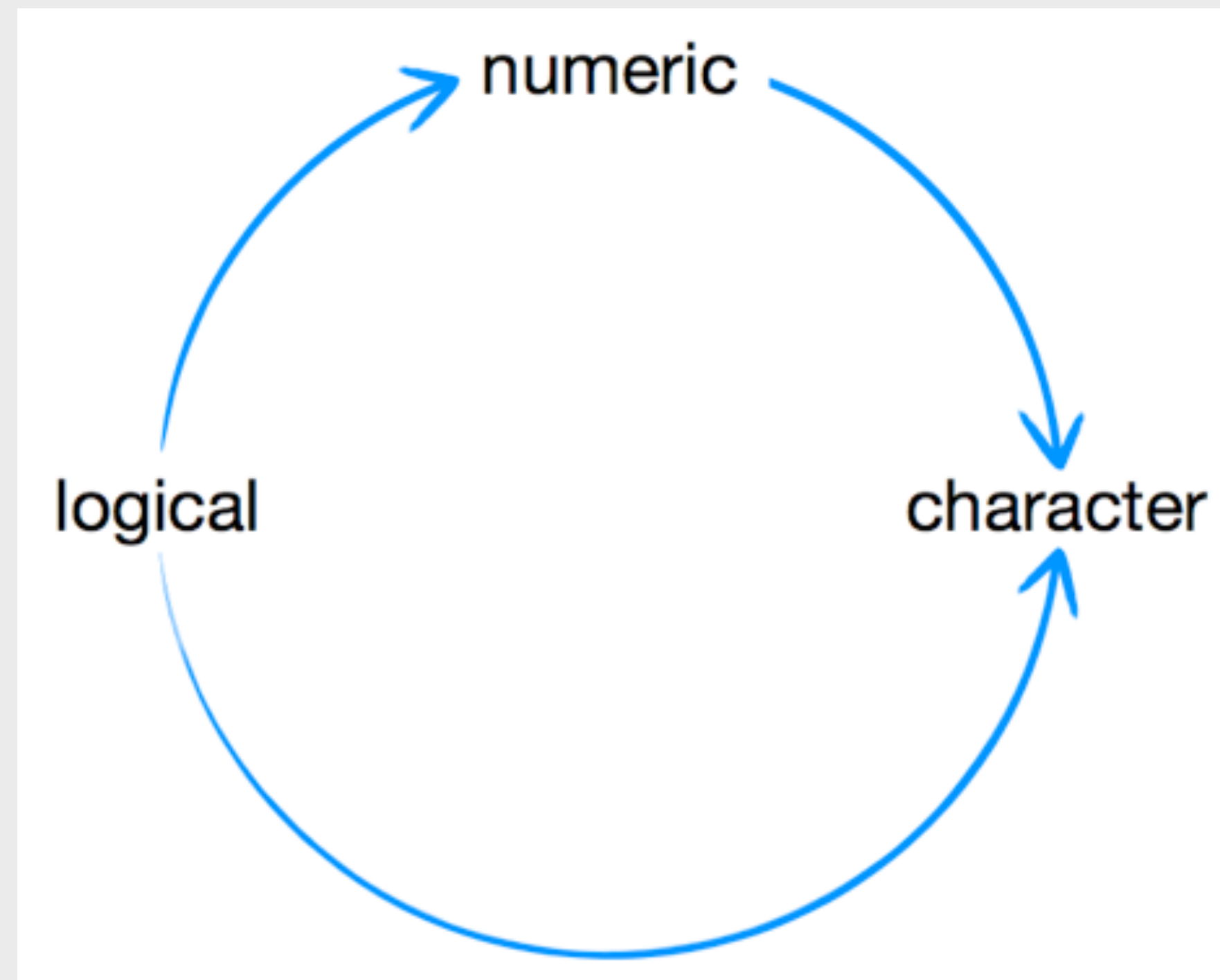
`c(TRUE, "a")`

`c(1, "TRUE")`

`TRUE + 5`

# Quiz

What type of data will result?



`c(5, "two")`  
character

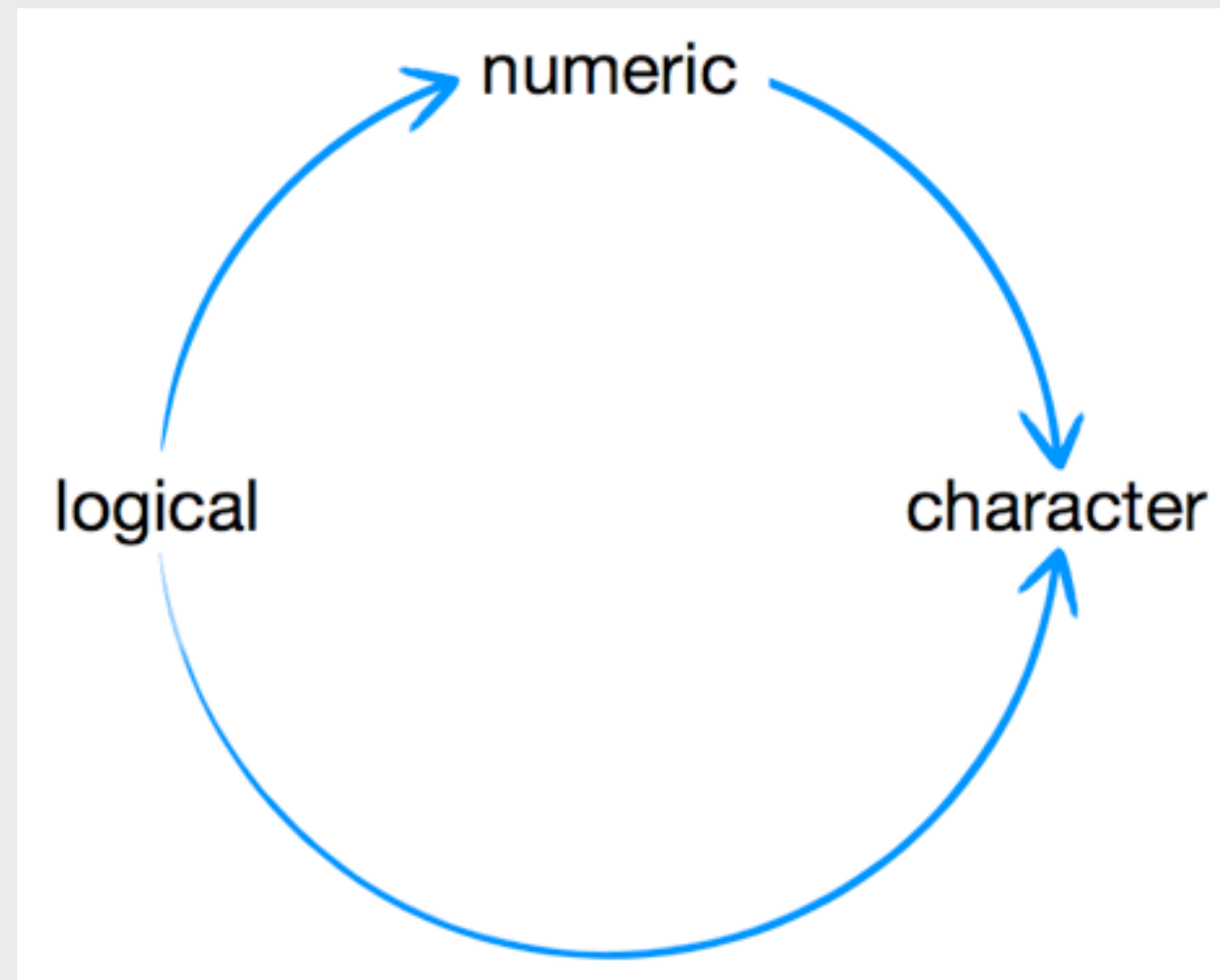
`c(TRUE, "a")`

`c(1, "TRUE")`

`TRUE + 5`

# Quiz

What type of data will result?



`c(5, "two")`  
character

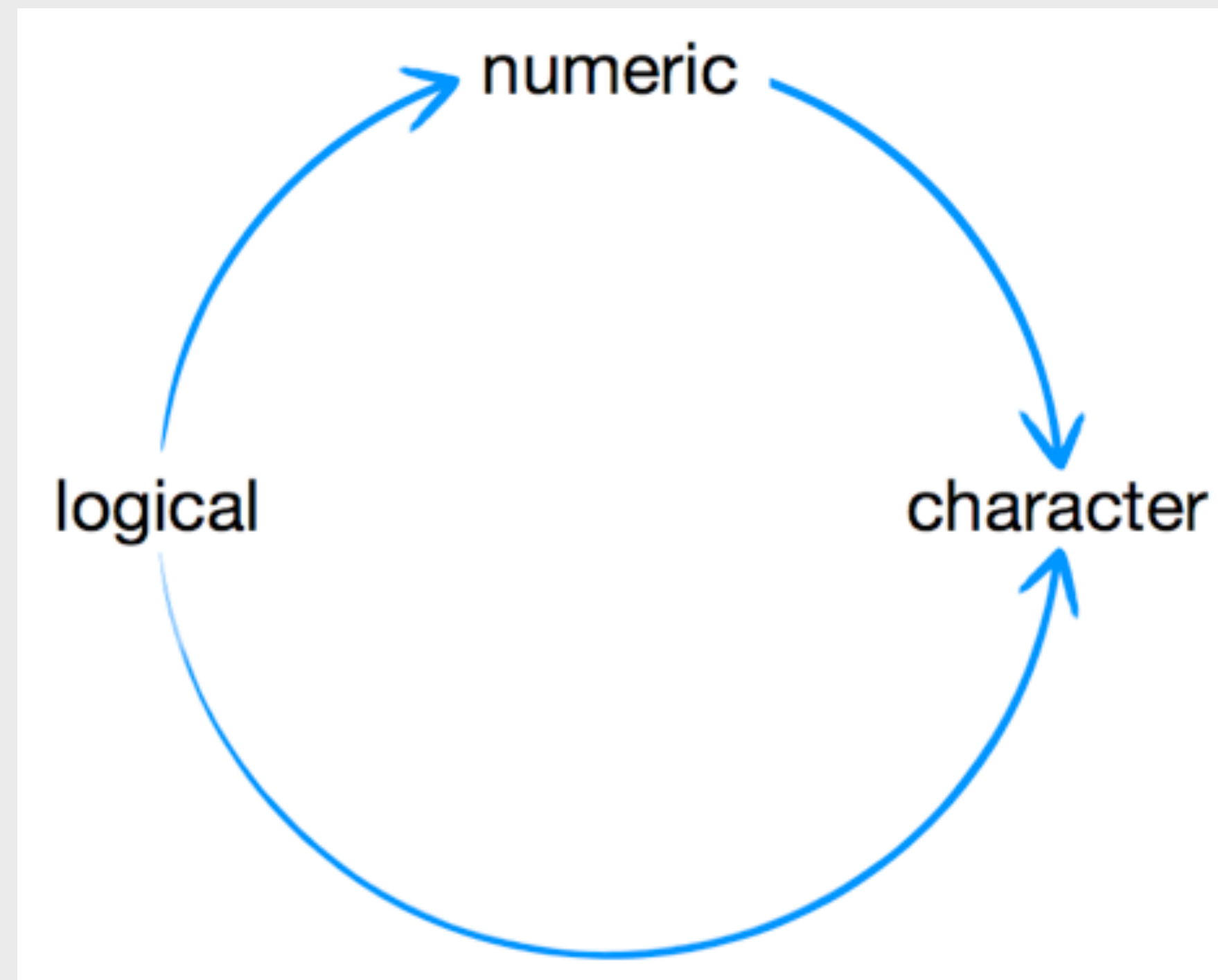
`c(TRUE, "a")`  
character

`c(1, "TRUE")`

`TRUE + 5`

# Quiz

What type of data will result?



`c(5, "two")`  
character

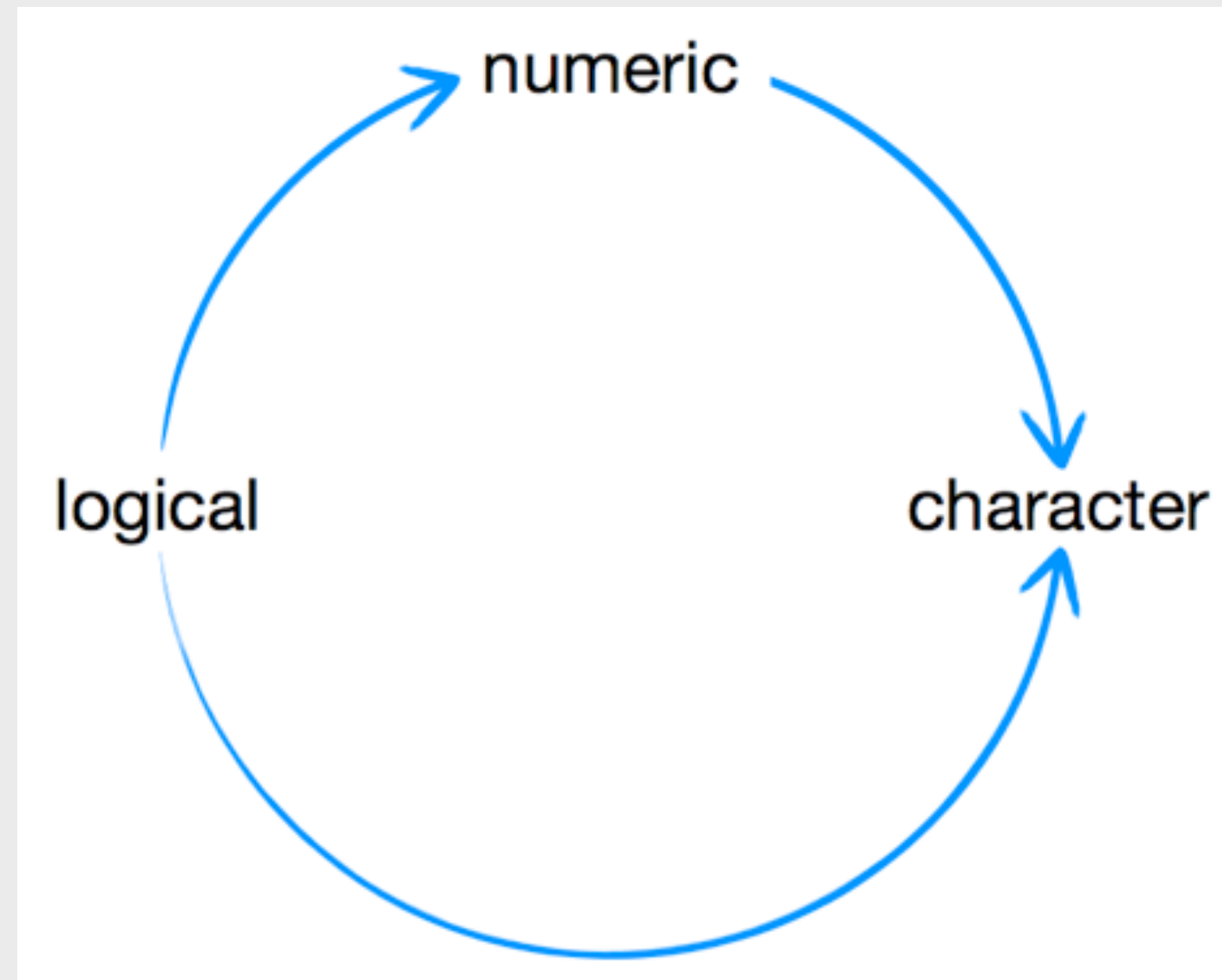
`c(TRUE, "a")`  
character

`c(1, "TRUE")`  
character

`TRUE + 5`

# Quiz

What type of data will result?



`c(5, "two")`  
character

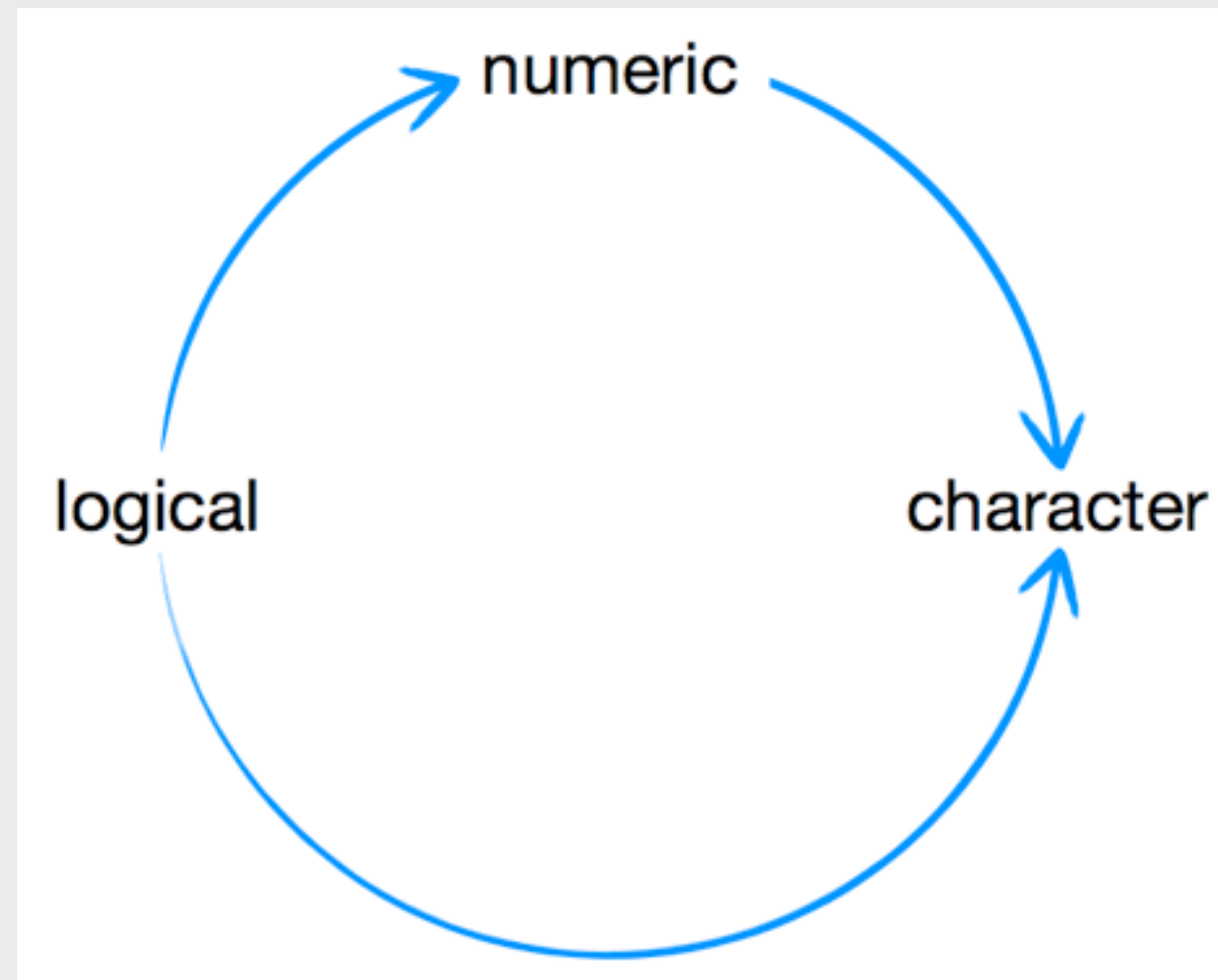
`c(TRUE, "a")`  
character

`c(1, "TRUE")`  
character

`TRUE + 5`

# Quiz

What type of data will result?



`c(5, "two")`  
character

`c(TRUE, "a")`  
character

`c(1, "TRUE")`  
character

`TRUE + 5`  
numeric

# manual coercion

function	coerces data to
as.numeric	numeric
as.character	character
as.logical	logical
as.factor	factor

`as.numeric("1")`

`as.character(TRUE)`

Matrix

1	"R"	TRUE
2	"S"	FALSE
3	"T"	TRUE

?



Matrix

"1"	"R"	"TRUE"
"2"	"S"	"FALSE"
"3"	"T"	"TRUE"

character

Matrix

"1"	"R"	"TRUE"
"2"	"S"	"FALSE"
"3"	"T"	"TRUE"

What if you want different data types in the same object?

# Lists and data frames

# lists and data frames

*lists* and *data frames* generalize vectors and matrices to allow multiple types of data

# Lists

# lists

A list is a one dimensional group of R objects.

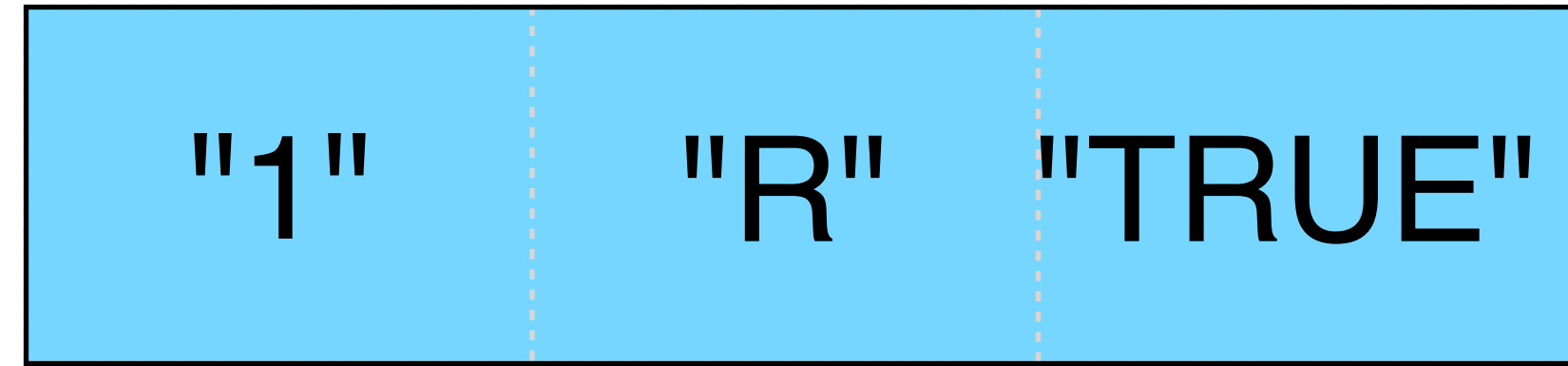
Create lists with `list`

```
lst <- list(1, "R", TRUE)
```

```
class(lst)
```

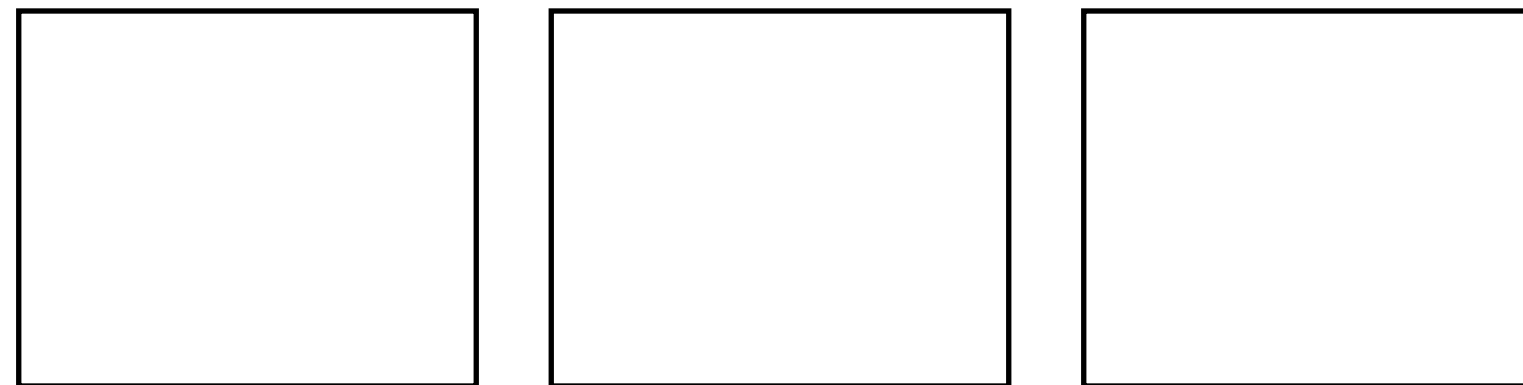
```
# "list"
```

Vector

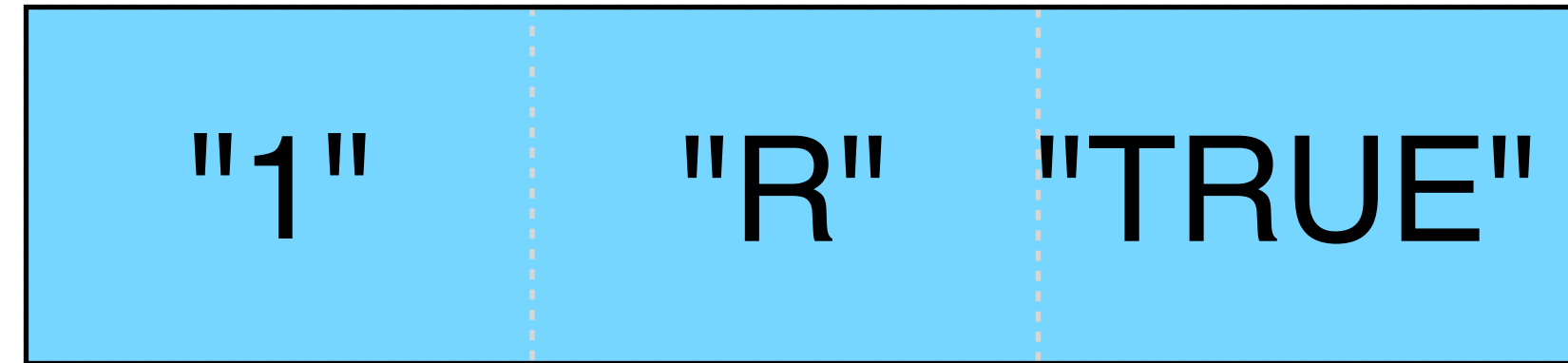


character

List



Vector



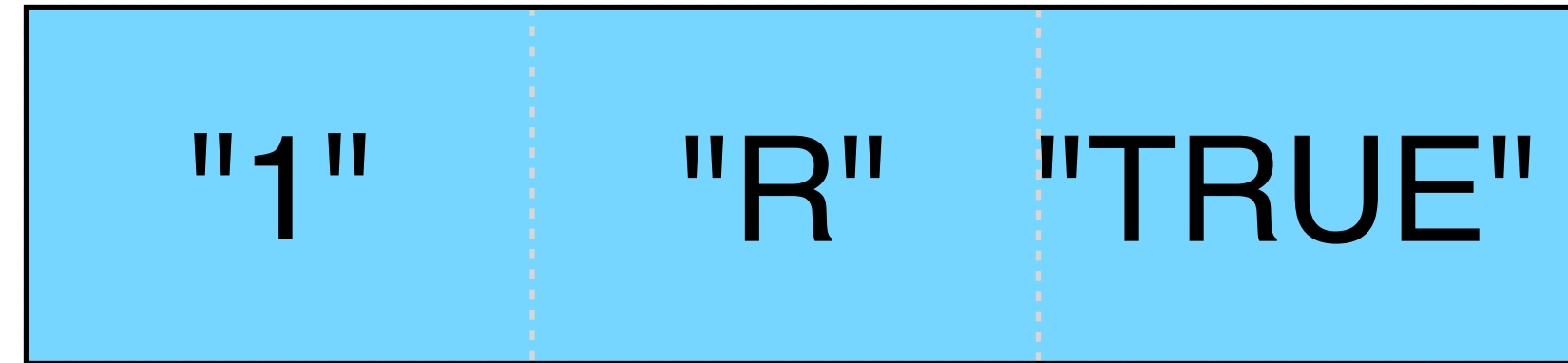
character

List





Vector



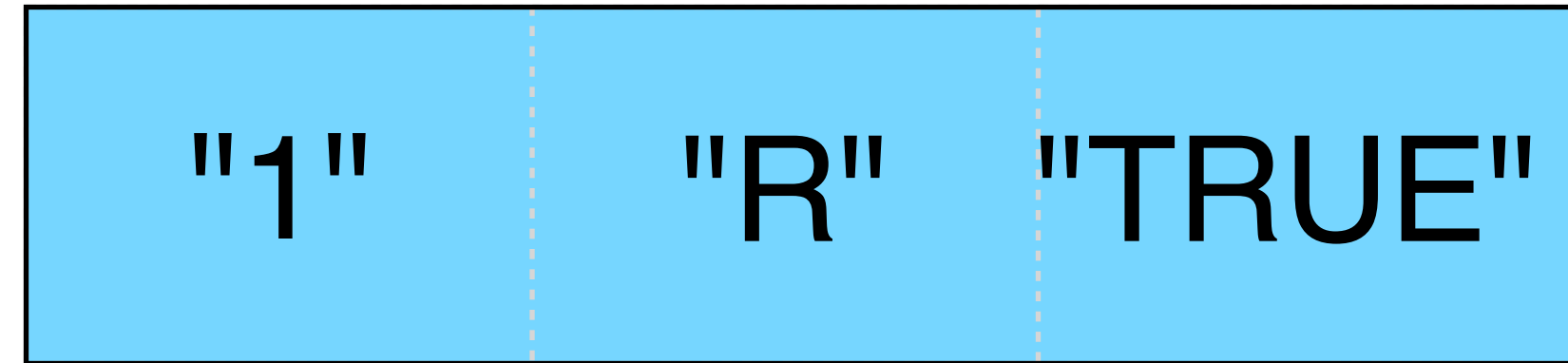
character

List



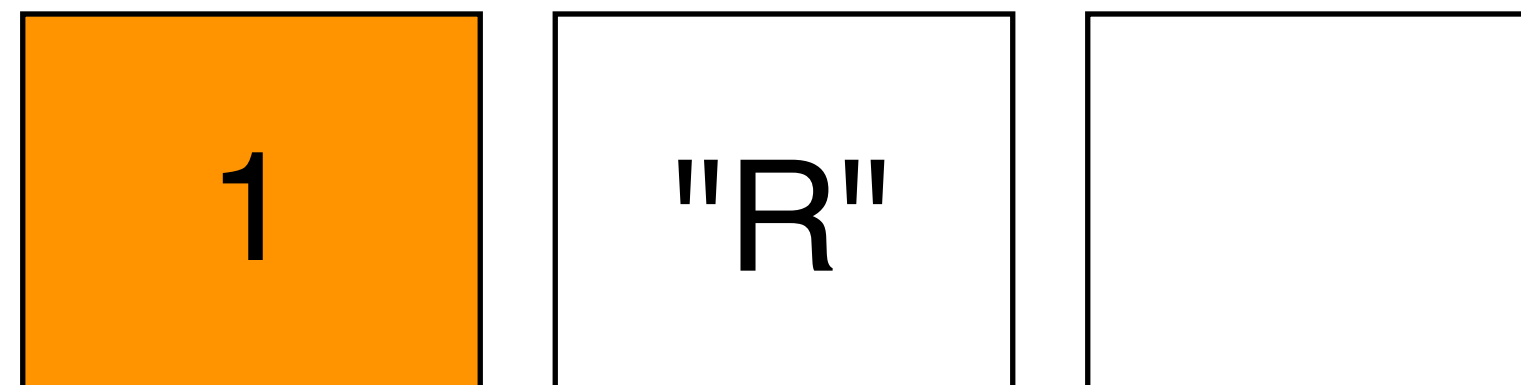
numeric

Vector



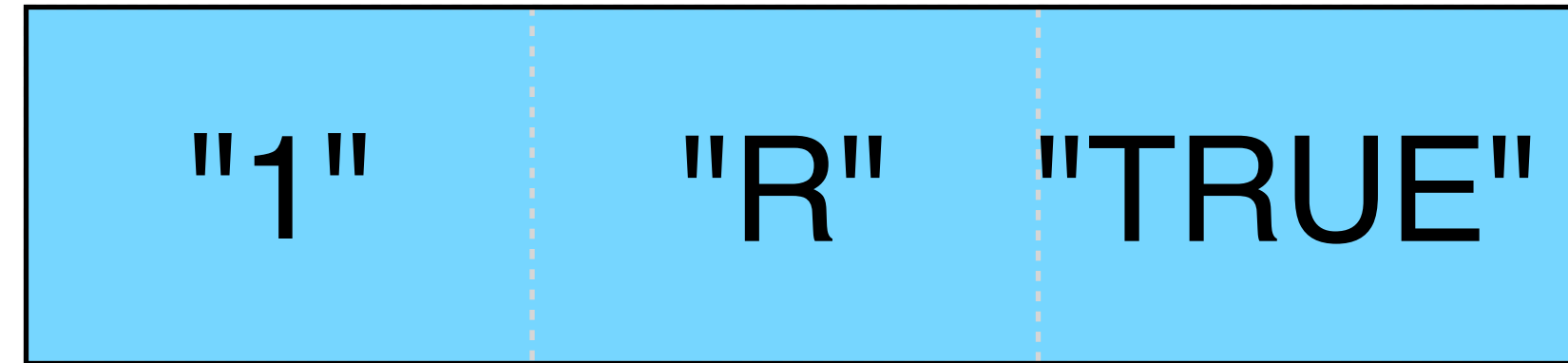
character

List



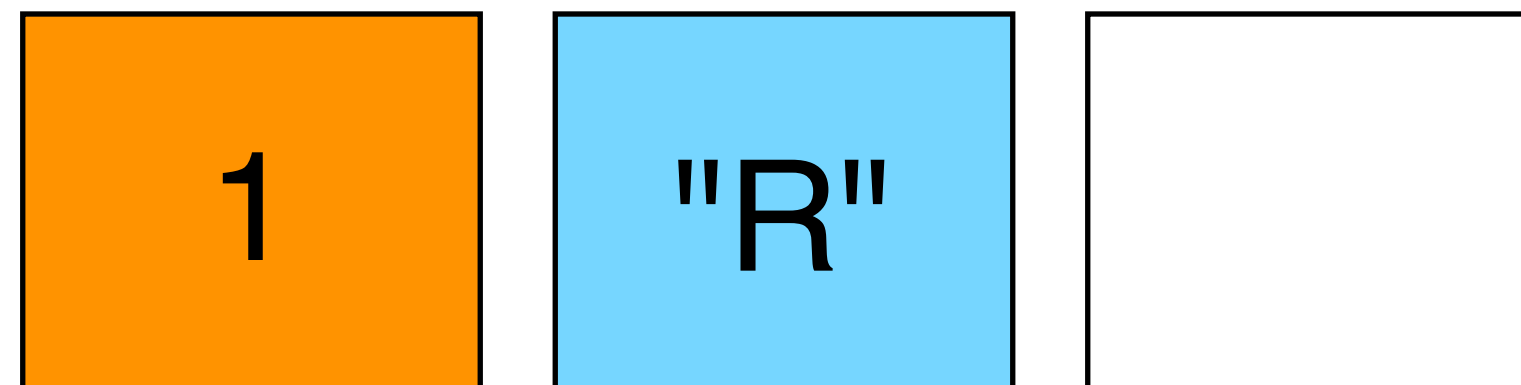
numeric

Vector



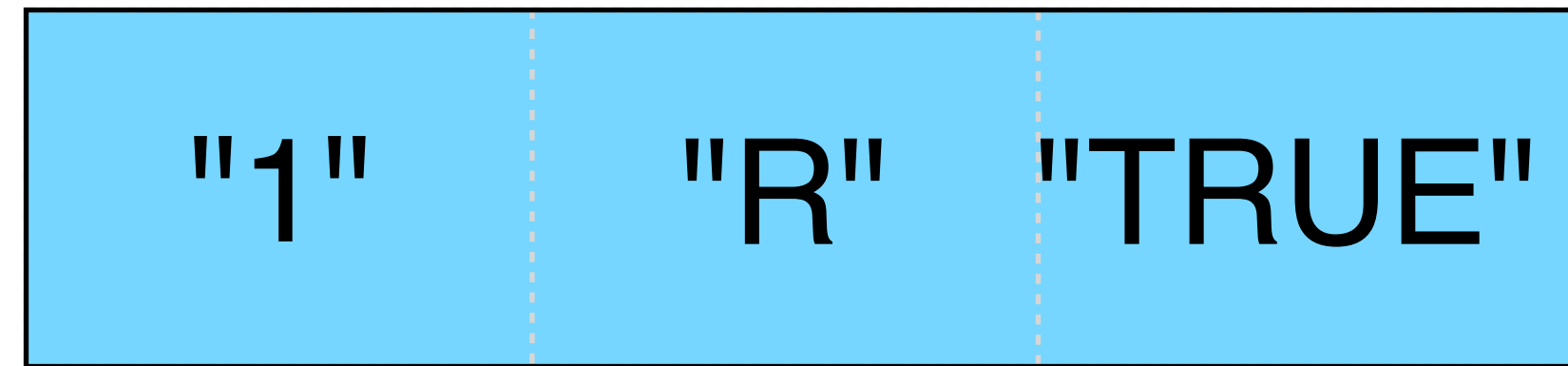
character

List



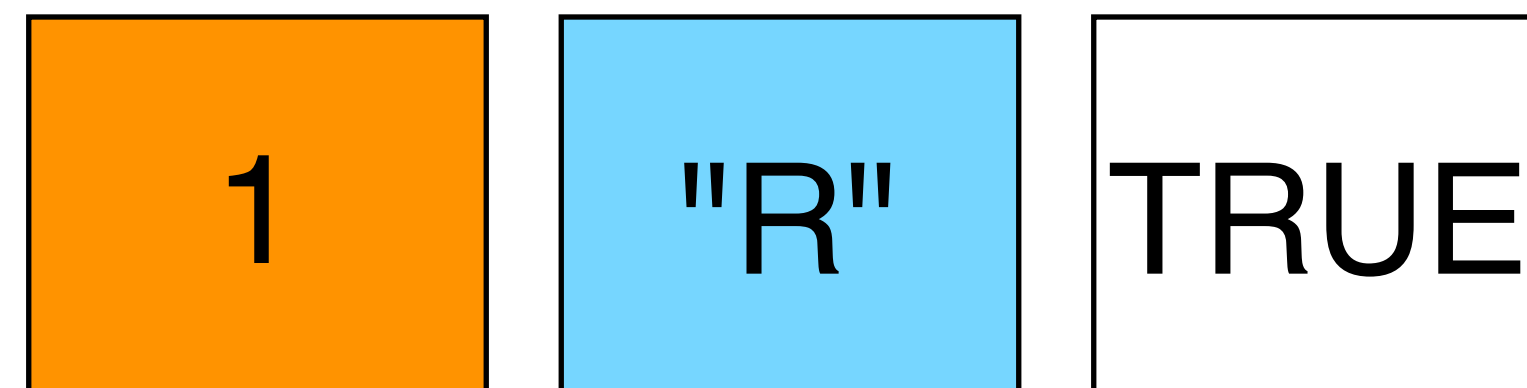
numeric character

Vector



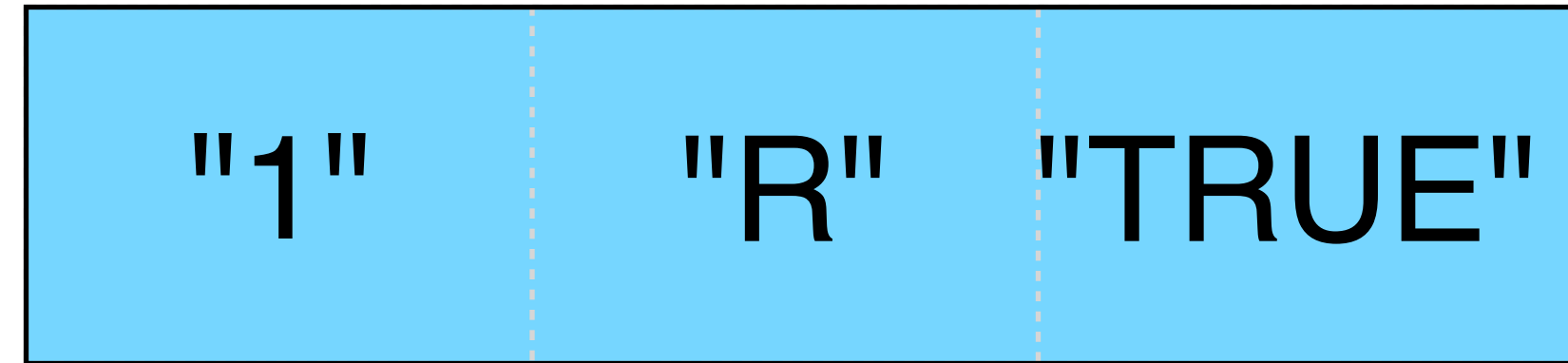
character

List



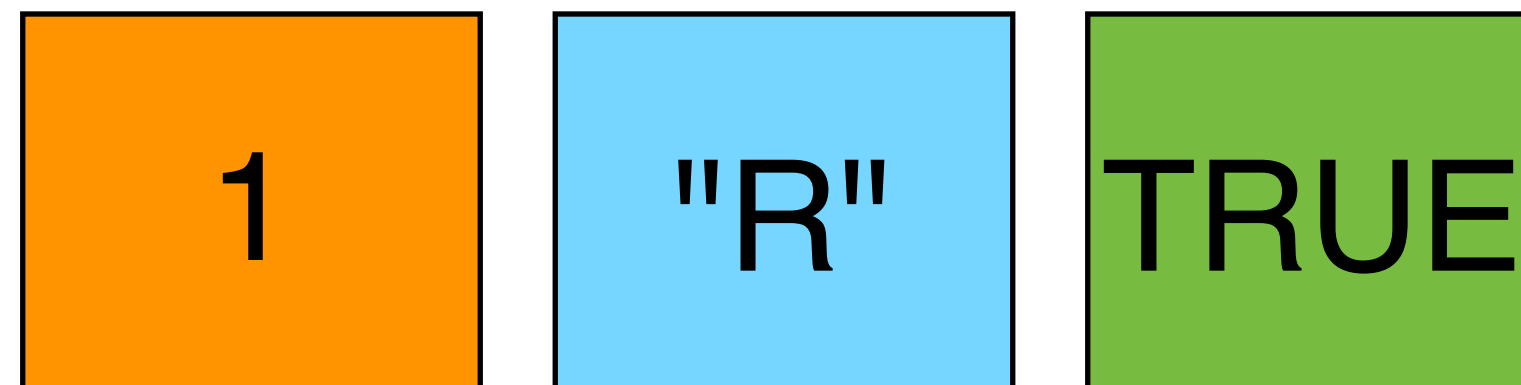
numeric character

Vector



character

List



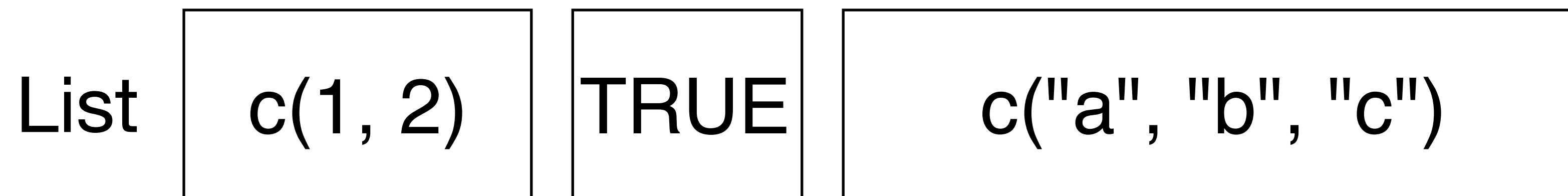
numeric

character

logical

The elements of a list can be anything. Even vectors or other lists.

```
list(c(1, 2), TRUE, c("a", "b", "c"))
```



# List viewer in RStudio

The screenshot shows the RStudio interface with the following components:

- Source Editor:** Contains R code for creating lists and data frames. The code is as follows:

```
94- ## Quiz
95- What type of data will result?
96-
97- ## Lists
98-
99- ```{r}
100- lst <- list(1, "R", TRUE)
101- class(lst)
102- ```
103-
104- ## Data frames
105- ```{r}
106- df <- data.frame(c(1, 2, 3),
107-   c("R","S","T"), c(TRUE, FALSE, TRUE))
108- class(df)
109- ```
110-
111- ## Naming
112- ```{r}
113- nvec <- c(one = 1, two = 2, three = 3)
114- nlst <- list(one = 1, two = 2, many = c(3, 4, 5))
115- ```
```
- Environment Pane:** Shows the Global Environment with the following objects:
  - Data:**
    - `bechdel`: 1794 obs. of 15 variables
    - `df`: 3 obs. of 3 variables
    - `lst`: List of 3
      - : num 1
      - : chr "R"
      - : logi TRUE
    - `mat`: num [1:2, 1:3] 1 2 3 4 5 6
    - `nlst`: List of 3
  - Values:**
    - `nvec`: Named num [1:3] 1 2 3
    - `threemill...`: 3e+06
    - `vec`: chr [1:3] "1" "R" "TRUE"
- Files Pane:** Shows the file explorer with the following files:
  - `02-Visualization.Rmd`: 2.5 KB, Jan 30, 2018, 5:33 PM
  - `01-Structures.Rmd`: 1.6 KB, Jan 30, 2018, 5:11 PM
  - `.Rhistory`: 131 B, Jan 28, 2018, 4:36 PM
  - `02-Visualization-solutions.Rmd`: 2.9 KB, Jan 30, 2018, 5:11 PM
  - `02-Visualization.nb.html`: 829.7 KB, Jan 30, 2018, 5:33 PM
  - `03-Syntax.Rmd`: 1.3 KB, Jan 30, 2018, 6:56 PM

Data frames



# data frame

A data frame is a two dimensional group of R objects.

Each column in a data frame can be a different type

```
df <- data.frame(c(1, 2, 3),  
  c("R", "S", "T"), c(TRUE, FALSE, TRUE))  
class(df)  
# "data.frame"
```

# Your turn

We've already seen a data frame today. What was it called? What kinds of data were in it?

00:30

Matrix

"1"	"R"	"TRUE"
"2"	"S"	"FALSE"
"3"	"T"	"TRUE"

character

data frame

1		
2		
3		

Matrix

"1"	"R"	"TRUE"
"2"	"S"	"FALSE"
"3"	"T"	"TRUE"

character

data frame

1		
2		
3		

numeric

Matrix

"1"	"R"	"TRUE"
"2"	"S"	"FALSE"
"3"	"T"	"TRUE"

character

data frame

1	"R"	
2	"S"	
3	"T"	

numeric

Matrix

"1"	"R"	"TRUE"
"2"	"S"	"FALSE"
"3"	"T"	"TRUE"

character

data frame

1	"R"	
2	"S"	
3	"T"	

numeric

character

Matrix

"1"	"R"	"TRUE"
"2"	"S"	"FALSE"
"3"	"T"	"TRUE"

character

data frame

1	"R"	TRUE
2	"S"	FALSE
3	"T"	TRUE

numeric

character

Matrix

"1"	"R"	"TRUE"
"2"	"S"	"FALSE"
"3"	"T"	"TRUE"

character

data frame

1	"R"	TRUE
2	"S"	FALSE
3	"T"	TRUE

numeric

character

logical



# names

You can name the elements of a vector, list, or data frame when you create them.

```
nvec <- c(one = 1, two = 2, three = 3)
```

```
nvec
```

```
# one two three
```

```
# 1 2 3
```

```
nlst <- list(one = 1, two = 2,  
            many = c(3, 4, 5))
```

```
nlst  
# $one  
# [1] 1  
#  
# $two  
# [1] 2  
#  
# $many  
# [1] 3 4 5
```

```
ndf <- data.frame(numbers = c(1, 2, 3),  
                  letters = c("R", "S", "T"),  
                  logic = c(TRUE, FALSE, TRUE))
```

```
ndf  
#  numbers letters logic  
# 1      1      R TRUE  
# 2      2      S FALSE  
# 3      3      T TRUE
```

# Your turn

Use the RStudio data preview to  
compare df and ndf

00:30

RStudio

Go to file/function | Addins

bechdel x 02-Visualization.Rmd x 03-Syntax.Rmd\* x ndf x df x

Filter

	numbers	letters	logic
1	1	R	TRUE
2	2	S	FALSE
3	3	T	TRUE

Showing 1 to 3 of 3 entries

Environment

Global Environment

Data

- bechdel
- df
- lst
- mat
- ndf
- nlst

Values

- nvec
- threem
- vec

Console Terminal x R Markdown x

~/Dropbox/Intro\_to\_R\_and\_RStudio/Day1/code/

```
> df <- data.frame(c(1, 2, 3),
```

Files Plots

New Folder

You can also see the names with names

```
names(ndf)
```

```
# [1] "numbers" "letters" "logic"
```

```
names(nvec)
```

```
# [1] "one" "two" "three"
```

single type

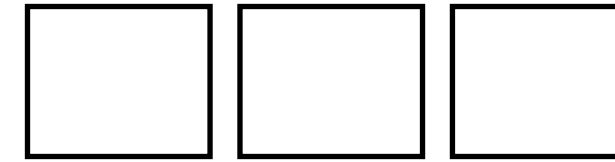
multiple types

1D

**Vector**

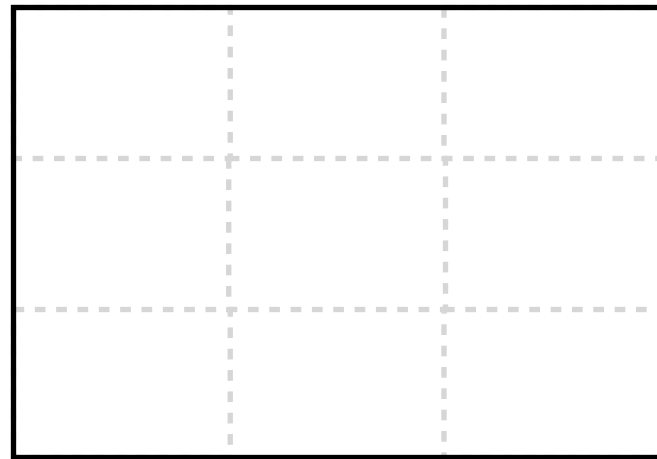


**List**

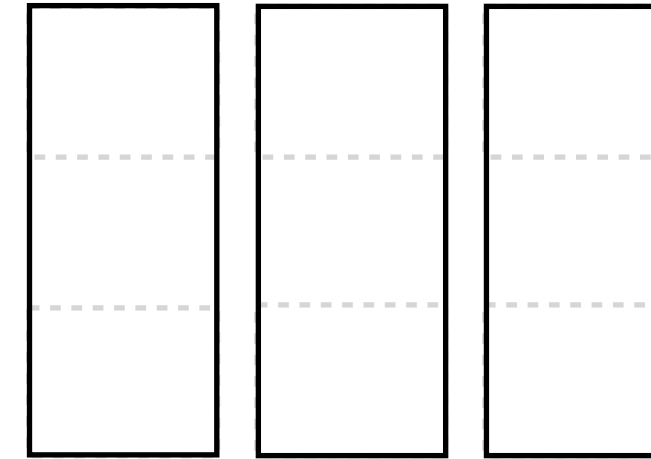


2D

**Matrix**

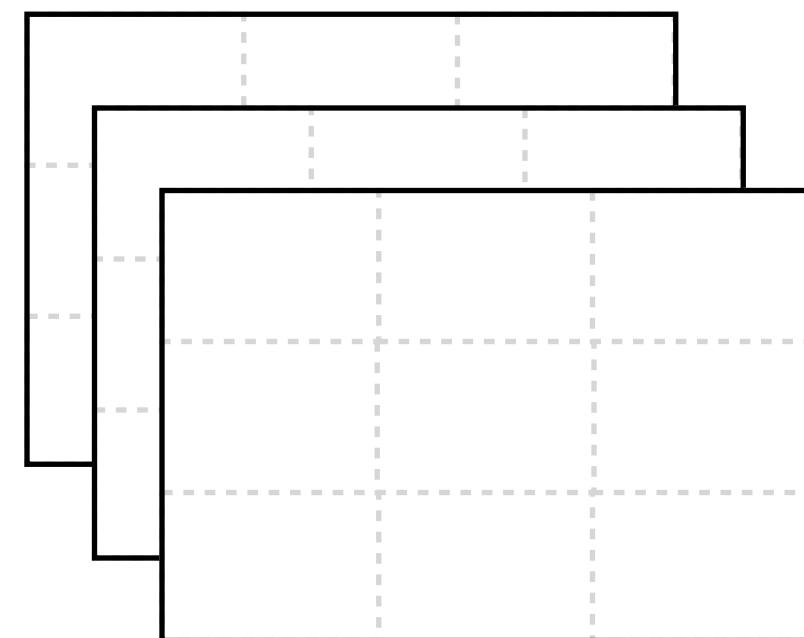


**Data frame**



nD

**Array**



# How R makes a data frame

List

```
c("a","b","c","d")
```

```
c(1, 2, 3, 4)
```

```
c(T, F, T, F)
```



List

```
c(  
  "a", "  
  b", "c  
  ", "d")
```

```
c(  
  1,  
  2,  
  3,  
  4)
```

```
c(  
  T,  
  F,  
  T,  
  F)
```

List

```
c(  
  "a", "  
  b", "c  
  ", "d")
```

```
c(  
  1,  
  2,  
  3,  
  4)
```

```
c(  
  T,  
  F,  
  T,  
  F)
```

~~List~~

data frame

```
c(
  "a", "
b", "c
", "d")
```

```
c(
  1,
  2,
  3,
  4)
```

```
c(
  T,
  F,
  T,
  F)
```

# helper functions for data structures

	create	change to	check	get names	get dimensions
vector	c, vector	as.vector	is.vector	names	length
matrix	matrix	as.matrix	is.matrix	rownames, colnames	dim, nrow, ncol
array	array	as.array	is.array	dimnames	dim
list	list	as.list	is.list	names	length
data frame	data.frame	as.data.frame	is.data.frame	names	dim, nrow, ncol

Syntax

**Syntax** is the set of rules that govern what code works and doesn't work in a programming language. Most programming languages offer one standardized syntax, but R allows package developers to specify their own syntax. As a result, there is a large variety of (equally valid) R syntaxes.

# R Syntax Comparison :: CHEAT SHEET

## Dollar sign syntax

```
goal(data$x, data$y)
```

### SUMMARY STATISTICS:

one continuous variable:  
`mean(mtcars$mpg)`

one categorical variable:  
`table(mtcars$cyl)`

two categorical variables:  
`table(mtcars$cyl, mtcars$am)`

one continuous, one categorical:  
`mean(mtcars$mpg[mtcars$cyl==4])`  
`mean(mtcars$mpg[mtcars$cyl==6])`  
`mean(mtcars$mpg[mtcars$cyl==8])`

### PLOTTING:

one continuous variable:  
`hist(mtcars$disp)`

`boxplot(mtcars$disp)`

one categorical variable:  
`barplot(table(mtcars$cyl))`

two continuous variables:  
`plot(mtcars$disp, mtcars$mpg)`

two categorical variables:  
`mosaicplot(table(mtcars$am, mtcars$cyl))`

one continuous, one categorical:  
`histogram(mtcars$disp[mtcars$cyl==4])`  
`histogram(mtcars$disp[mtcars$cyl==6])`  
`histogram(mtcars$disp[mtcars$cyl==8])`

`boxplot(mtcars$disp[mtcars$cyl==4])`  
`boxplot(mtcars$disp[mtcars$cyl==6])`  
`boxplot(mtcars$disp[mtcars$cyl==8])`

### WRANGLING:

subsetting:  
`mtcars[mtcars$mpg>30, ]`

making a new variable:  
`mtcars$efficient[mtcars$mpg>30] <- TRUE`  
`mtcars$efficient[mtcars$mpg<30] <- FALSE`

## Formula syntax

```
goal(y~x|z, data=data, group=w)
```

### SUMMARY STATISTICS:

one continuous variable:  
`mosaic::mean(~mpg, data=mtcars)`

one categorical variable:  
`mosaic::tally(~cyl, data=mtcars)`

two categorical variables:  
`mosaic::tally(cyl~am, data=mtcars)`

one continuous, one categorical:  
`mosaic::mean(mpg~cyl, data=mtcars)`

tilde

### PLOTTING:

one continuous variable:  
`lattice::histogram(~disp, data=mtcars)`

`lattice::bwplot(~disp, data=mtcars)`

one categorical variable:  
`mosaic::bargraph(~cyl, data=mtcars)`

two continuous variables:  
`lattice::xyplot(mpg~disp, data=mtcars)`

two categorical variables:  
`mosaic::bargraph(~am, data=mtcars, group=cyl)`

one continuous, one categorical:  
`lattice::histogram(~disp|cyl, data=mtcars)`  
`lattice::bwplot(cyl~disp, data=mtcars)`

## Tidyverse syntax

```
data %>% goal(x)
```

### SUMMARY STATISTICS:

one continuous variable:  
`mtcars %>% dplyr::summarize(mean(mpg))`

one categorical variable:  
`mtcars %>% dplyr::group_by(cyl) %>%`  
`dplyr::summarize(n())`

two categorical variables:  
`mtcars %>% dplyr::group_by(cyl, am) %>%`  
`dplyr::summarize(n())`

one continuous, one categorical:  
`mtcars %>% dplyr::group_by(cyl) %>%`  
`dplyr::summarize(mean(mpg))`

the pipe

### PLOTTING:

one continuous variable:  
`ggplot2::qplot(x=mpg, data=mtcars, geom="histogram")`

`ggplot2::qplot(y=disp, x=1, data=mtcars, geom="boxplot")`

one categorical variable:  
`ggplot2::qplot(x=cyl, data=mtcars, geom="bar")`

two continuous variables:  
`ggplot2::qplot(x=disp, y=mpg, data=mtcars, geom="point")`

two categorical variables:  
`ggplot2::qplot(x=factor(cyl), data=mtcars, geom="bar") +`  
`facet_grid(~am)`

one continuous, one categorical:  
`ggplot2::qplot(x=disp, data=mtcars, geom="histogram") +`  
`facet_grid(~cyl)`

`ggplot2::qplot(y=disp, x=factor(cyl), data=mtcars,`  
`geom="boxplot")`

### WRANGLING:

subsetting:  
`mtcars %>% dplyr::filter(mpg>30)`

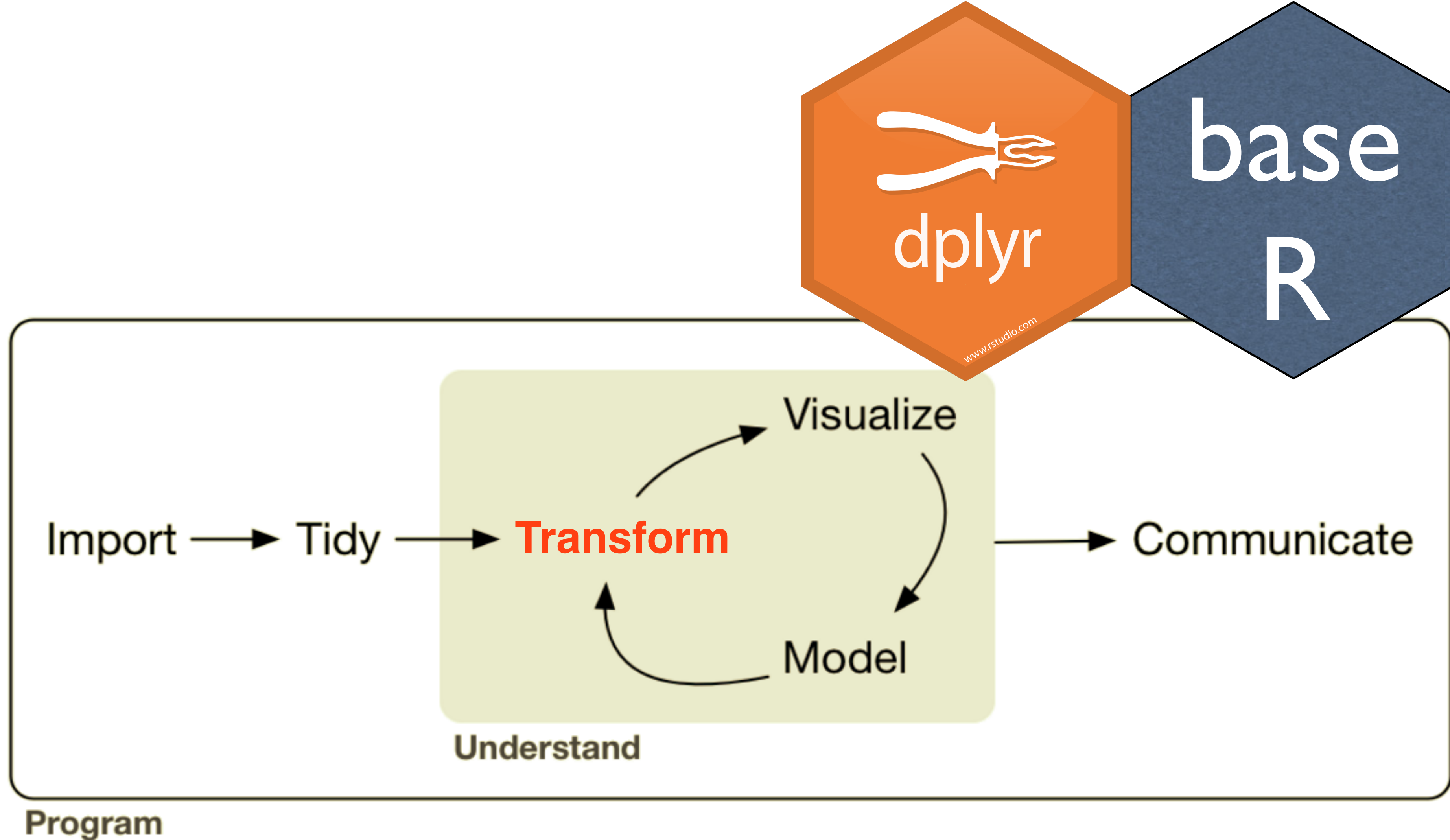
making a new variable:  
`mtcars <- mtcars %>%`  
`dplyr::mutate(efficient = if_else(mpg>30, TRUE, FALSE))`

The variety of R syntaxes give you many ways to “say” the same thing

read across the cheatsheet to see how different syntaxes approach the same problem

# Subsetting





From *R for Data Science* by Hadley Wickham and Garrett Grolemund.

# Toy data

```
beatles <- data.frame(  
  name = c("John", "Paul", "George", "Ringo"),  
  birth = c(1940, 1942, 1943, 1940),  
  instrument = c("guitar", "bass", "guitar", "drums")  
)
```

First— the  
tidyverse way:  
dplyr

# dplyr methods for isolating data

**select()** - extract **variables**  
**filter()** - extract **cases**  
**arrange()** - reorder **cases**



# select()

Extract columns by name.

```
select(.data, ...)
```

**data frame  
to  
transform**

**name(s) of columns to  
extract**  
(or a select helper function)



# select()

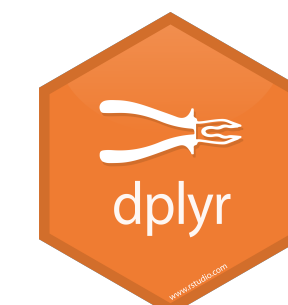
Extract columns by name.

```
select(beatles, name, birth)
```

name	birth	instrument
John	1940	guitar
Paul	1942	base
George	1943	guitar
Ringo	1940	drums

→

name	birth
John	1940
Paul	1942
George	1943
Ringo	1940





# Your Turn

Alter the code to select just the **instrument** column:

```
select(beatles, name, birth)
```

01:00



# select() helpers

**:** - Select range of columns

```
select(storms, storm:pressure)
```

**-** - Select every column but

```
select(storms, -c(storm, pressure))
```

**starts\_with()** - Select columns that start with...

```
select(storms, starts_with("w"))
```

**ends\_with()** - Select columns that end with...

```
select(storms, ends_with("e"))
```





# select() helpers

**contains()** - Select columns whose names contain...

```
select(storms, contains("d"))
```

**matches()** - Select columns whose names match regular expression

```
select(storms, matches("^.{4}$"))
```

**one\_of()** - Select columns whose names are one of a set

```
select(storms, one_of(c("storm", "storms", "Storm")))
```

**num\_range()** - Select columns named in prefix, number style

```
select(storms, num_range("x", 1:5))
```

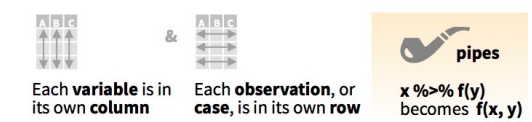


# select() helpers

## Data Transformation with dplyr : : CHEAT SHEET

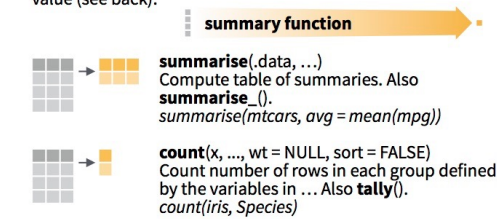


dplyr functions work with pipes and expect tidy data. In tidy data:



### Summarise Cases

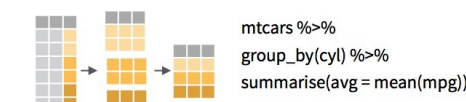
These apply **summary functions** to columns to create a new table. Summary functions take vectors as input and return one value (see back).



**VARIATIONS**  
**summarise\_all()** - Apply funs to every column.  
**summarise\_at()** - Apply funs to specific columns.  
**summarise\_if()** - Apply funs to all cols of one type.

### Group Cases

Use **group\_by()** to create a "grouped" copy of a table. dplyr functions will manipulate each "group" separately and then combine the results.



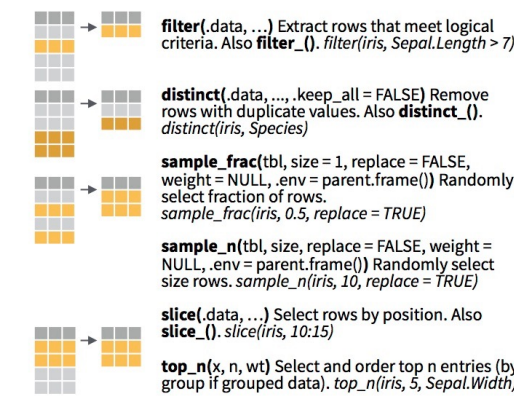
**group\_by()**(data, ..., add = FALSE)  
 Returns copy of table grouped by ...  
`g_iris <- group_by(iris, Species)`

**ungroup()**(x, ...)  
 Returns ungrouped copy of table.  
`ungroup(g_iris)`

### Manipulate Cases

#### EXTRACT CASES

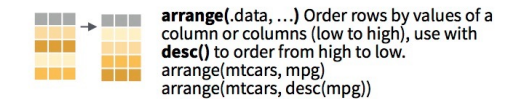
Row functions return a subset of rows as a new table. Use a variant that ends in `_` for non-standard evaluation friendly code.



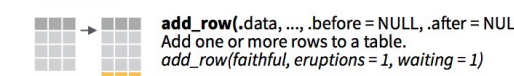
#### Logical and boolean operators to use with filter()

< is.na() %in% | xor()  
 <= is.na() ! %in% | &  
 > !is.na() !  
 See ?base::logic and ?Comparison for help.

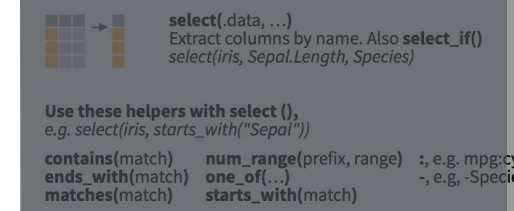
#### ARRANGE CASES



#### ADD CASES

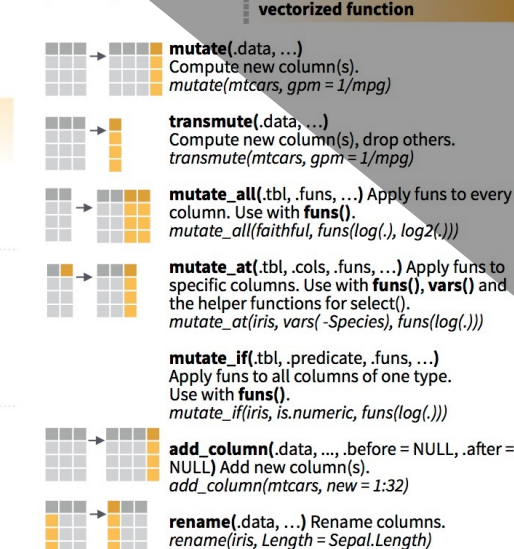


Column functions return a set of columns as a new table. Use a variant that ends in `_` for non-standard evaluation friendly code.



#### MAKE NEW VARIABLES

These apply **vectorized functions** to columns. Vectorized funs take vectors as input and return vectors of the same length as output (see back).



Column functions return a set of columns as a new table. Use a variant that ends in `_` for non-standard evaluation friendly code.



**select()**(data, ...)  
 Extract columns by name. Also **select\_if()**  
`select(iris, Sepal.Length, Species)`

Use these helpers with **select\_()**,  
 e.g. `select(iris, starts_with("Sepal"))`

**contains()**(match) `num_range(prefix, range)` ;, e.g. `mpg:cyl`  
**ends\_with()**(match) `one_of(...)` -, e.g. `-Species`  
**matches()**(match) `starts_with(match)`



Now, the base R  
way: brackets  
and dollar signs

# Base R bracket subset notation

in base R, you use the same syntax to

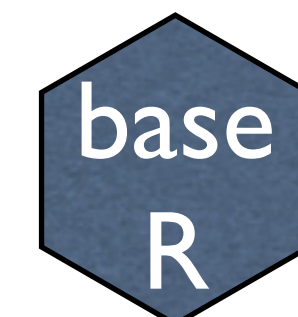
extract **variables**

extract **cases**

name of object  
to subset

brackets  
(brackets always mean  
subset)

`vec[ ]`



# Subset notation

name of object  
to subset

**vec**



# Subset notation

name of object  
to subset

brackets  
(brackets always mean  
subset)

`vec[?]`

an index  
(that tells R which  
elements to include)

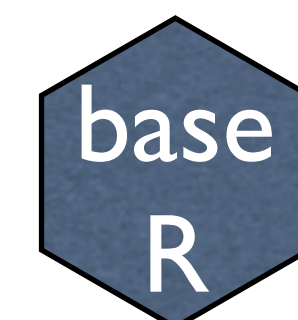




Each dimension needs its own index!

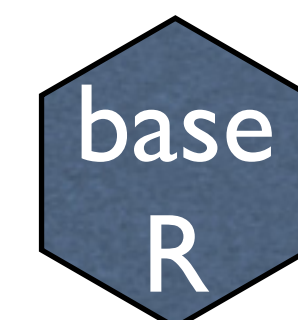
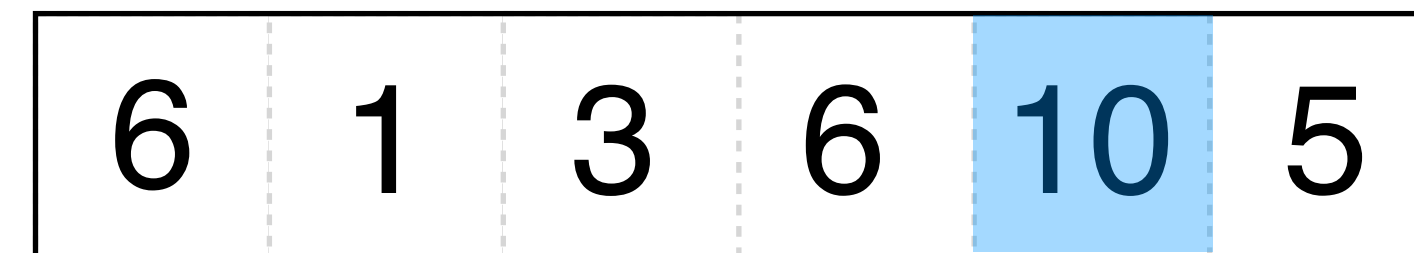
vec[?]

6	1	3	6	10	5
---	---	---	---	----	---



Each dimension needs its own index!

vec[?]





Each dimension needs its own index!

`vec[?]`  
`beatles[?, ?]`

John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

# Each dimension needs its own index!

`vec[?]`  
`beatles[?, ?]`

which  
**rows** to  
include

John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

# Each dimension needs its own index!

vec[?]  
beatles[?, ?]

John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

which  
**rows** to  
include

which  
**columns**  
to include

# Each dimension needs its own index!

`vec[?]`  
`beatles[?, ?]`

John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

which  
**rows** to  
include

separate  
dimensions  
with a  
**comma**

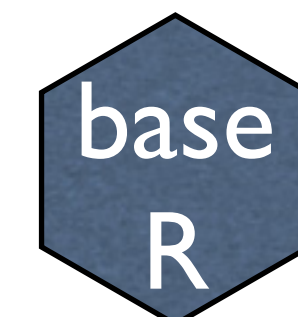
which  
**columns**  
to include

Each dimension needs its own index!

vec[?]  
beatles[?, ?]

John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

What goes in the indexes?



# Four ways to subset

1. Integers

2. Blank spaces

3. Names

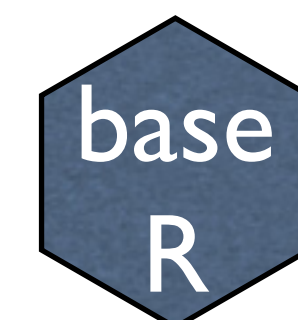
4. Logical vectors (TRUE and FALSE)

# Integers (positive)

Positive integers behave just like  $ij$  notation in linear algebra

beatles[?, ?]


John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums



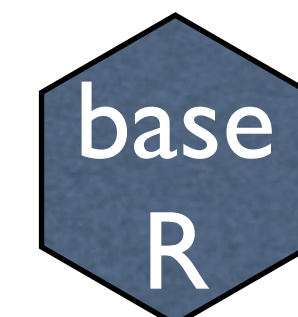
# Integers (positive)

Positive integers behave just like  $ij$  notation in linear algebra

beatles[2, ?]



John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

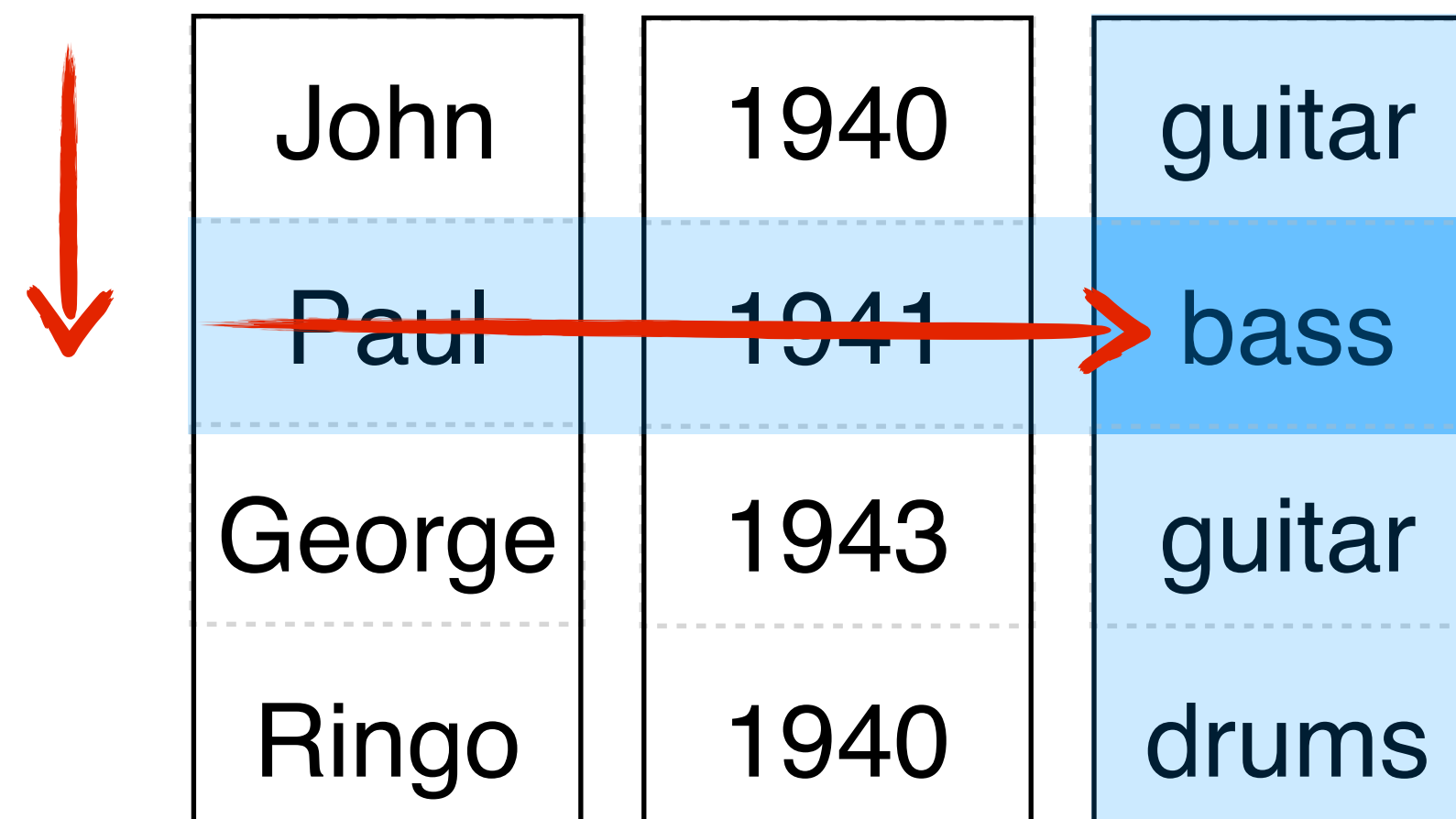




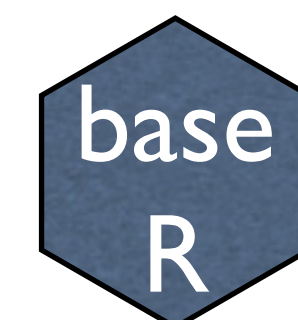
# Integers (positive)

Positive integers behave just like  $ij$  notation in linear algebra

beatles[2,3]



John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

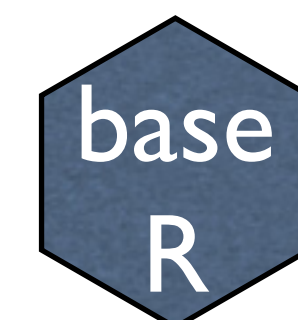


# Integers (positive)

Positive integers behave just like *ij* notation in linear algebra

beatles[2,3]

John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums



```
c("John", "Paul",  
"George", "Ringo")
```

```
c(1940, 1942,  
1943, 1940)
```

```
c("guitar",  
"drums")
```

John	1940	
Paul	1942	
George	1943	guitar
Ringo	1940	drums

`beatles[2]`

Numeric indexing is a sign of old R code, and there are few good use cases  
We won't focus on it

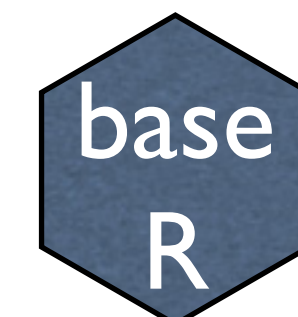


# Names

If your object has names, you can ask for elements or columns back by name.

```
beatles[, "birth"]
```

name	birth	instrument
John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

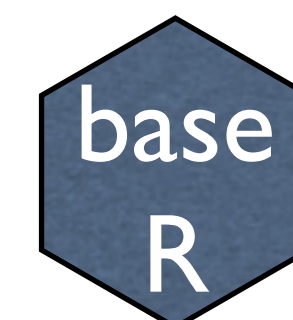


# Names

If your object has names, you can ask for elements or columns back by name.

```
beatles[,c("name","birth")]
```

name	birth	instrument
John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums





# Your Turn

Modify the code below to select just the `instrument` column

```
beatles[ , "birth"]
```

00:30



# \$

The most common syntax for subsetting lists  
and data frames



name	birth	instrument
John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

beatles\$birth





name	birth	instrument
John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

`beatles$birth`

name of  
data frame



name	birth	instrument
John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

`beatles$birth`

name of  
data frame

\$



name	birth	instrument
John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

`beatles$birth`

name of  
data frame

\$

name of column  
(no quotes)



name	birth	instrument
John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

`c(1940, 1941, 1943, 1940)`

`beatles$birth`

name of  
data frame

\$

name of column  
(no quotes)





# Your Turn

Modify the code below to select just the  
`instrument` column

```
beatles$birth
```

00:30



# Logical comparisons

# Logical comparisons

?Comparison

<code>x &lt; y</code>	Less than
<code>x &gt; y</code>	Greater than
<code>x == y</code>	Equal to
<code>x &lt;= y</code>	Less than or equal to
<code>x &gt;= y</code>	Greater than or equal to
<code>x != y</code>	Not equal to
<code>x %in% y</code>	Group membership
<code>is.na(x)</code>	Is NA
<code>!is.na(x)</code>	Is not NA





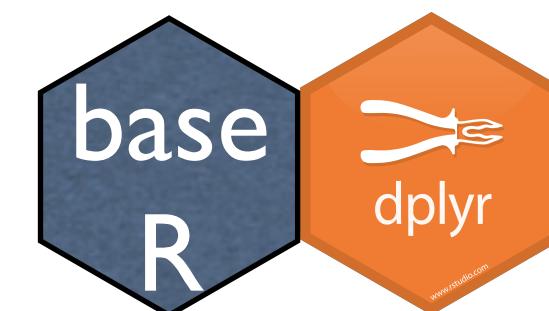
# Logical comparisons

What will these return?

`1 < 3`

`1 > 3`

`c(1, 2, 3, 4, 5) > 3`



# %in%

# What does this do?

```
1 %in% c(1, 2, 3, 4)
```

```
1 %in% c(2, 3, 4)
```

```
c(3,4,5,6) %in% c(2, 3, 4)
```



# %in%

%in% tests whether the object on the left is a member of the group on the right.

```
1 %in% c(1, 2, 3, 4)
```

```
# TRUE
```

```
1 %in% c(2, 3, 4)
```

```
# FALSE
```

```
c(3,4,5,6) %in% c(2, 3, 4)
```

```
# TRUE TRUE FALSE FALSE
```



# Your turn

`x <- c(1, 2, 3, 4, 5)`

Operator	Result	Comparison
<code>x &gt; 3</code>	<code>c(F, F, F, T, T)</code>	greater than
<code>x &gt;= 3</code>		
<code>x &lt; 3</code>		
<code>x &lt;= 3</code>		
<code>x == 3</code>		
<code>x != 3</code>		
<code>x = 3</code>		

01:00

# Your turn

`x <- c(1, 2, 3, 4, 5)`

Operator	Result	Comparison
<code>x &gt; 3</code>	<code>c(F, F, F, T, T)</code>	greater than
<code>x &gt;= 3</code>	<code>c(F, F, T, T, T)</code>	greater than or equal to
<code>x &lt; 3</code>		
<code>x &lt;= 3</code>		
<code>x == 3</code>		
<code>x != 3</code>		
<code>x = 3</code>		

# Your turn

`x <- c(1, 2, 3, 4, 5)`

Operator	Result	Comparison
<code>x &gt; 3</code>	<code>c(F, F, F, T, T)</code>	greater than
<code>x &gt;= 3</code>	<code>c(F, F, T, T, T)</code>	greater than or equal to
<code>x &lt; 3</code>	<code>c(T, T, F, F, F)</code>	less than
<code>x &lt;= 3</code>		
<code>x == 3</code>		
<code>x != 3</code>		
<code>x = 3</code>		

# Your turn

`x <- c(1, 2, 3, 4, 5)`

Operator	Result	Comparison
<code>x &gt; 3</code>	<code>c(F, F, F, T, T)</code>	greater than
<code>x &gt;= 3</code>	<code>c(F, F, T, T, T)</code>	greater than or equal to
<code>x &lt; 3</code>	<code>c(T, T, F, F, F)</code>	less than
<code>x &lt;= 3</code>	<code>c(T, T, T, F, F)</code>	less than or equal to
<code>x == 3</code>		
<code>x != 3</code>		
<code>x = 3</code>		



# Your turn

`x <- c(1, 2, 3, 4, 5)`

Operator	Result	Comparison
<code>x &gt; 3</code>	<code>c(F, F, F, T, T)</code>	greater than
<code>x &gt;= 3</code>	<code>c(F, F, T, T, T)</code>	greater than or equal to
<code>x &lt; 3</code>	<code>c(T, T, F, F, F)</code>	less than
<code>x &lt;= 3</code>	<code>c(T, T, T, F, F)</code>	less than or equal to
<code>x == 3</code>	<code>c(F, F, T, F, F)</code>	equal to
<code>x != 3</code>		
<code>x = 3</code>		

# Your turn

`x <- c(1, 2, 3, 4, 5)`

Operator	Result	Comparison
<code>x &gt; 3</code>	<code>c(F, F, F, T, T)</code>	greater than
<code>x &gt;= 3</code>	<code>c(F, F, T, T, T)</code>	greater than or equal to
<code>x &lt; 3</code>	<code>c(T, T, F, F, F)</code>	less than
<code>x &lt;= 3</code>	<code>c(T, T, T, F, F)</code>	less than or equal to
<code>x == 3</code>	<code>c(F, F, T, F, F)</code>	equal to
<code>x != 3</code>	<code>c(T, T, F, T, T)</code>	not equal to
<code>x = 3</code>		

# Your turn

`x <- c(1, 2, 3, 4, 5)`

Operator	Result	Comparison
<code>x &gt; 3</code>	<code>c(F, F, F, T, T)</code>	greater than
<code>x &gt;= 3</code>	<code>c(F, F, T, T, T)</code>	greater than or equal to
<code>x &lt; 3</code>	<code>c(T, T, F, F, F)</code>	less than
<code>x &lt;= 3</code>	<code>c(T, T, T, F, F)</code>	less than or equal to
<code>x == 3</code>	<code>c(F, F, T, F, F)</code>	equal to
<code>x != 3</code>	<code>c(T, T, F, T, T)</code>	not equal to
<code>x = 3</code>		same as <-

# Boolean operators

# Boolean operators

?base::Logic

<code>a &amp; b</code>	and
<code>a   b</code>	or
<code>xor(a, b)</code>	exactly or
<code>!a</code>	not



# Boolean operators

You can combine logical tests with `&`, `|`, `xor`, `!`, `any`, and `all`

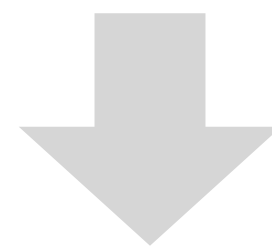
`x > 2 & x < 9`



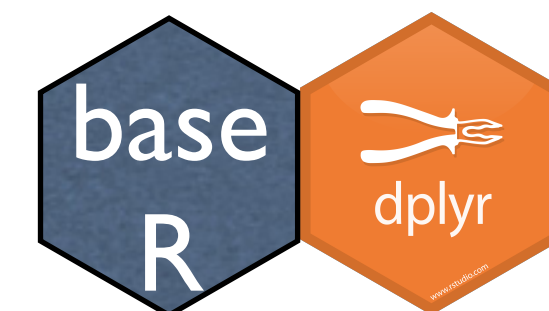
# Boolean operators

You can combine logical tests with `&`, `|`, `xor`, `!`, `any`, and `all`

`x > 2 & x < 9`



`TRUE &`

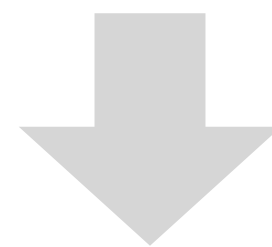




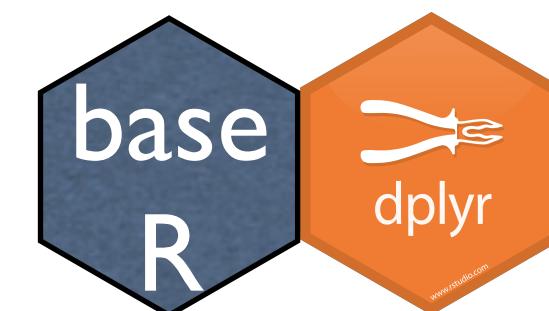
# Boolean operators

You can combine logical tests with `&`, `|`, `xor`, `!`, `any`, and `all`

`x > 2 & x < 9`



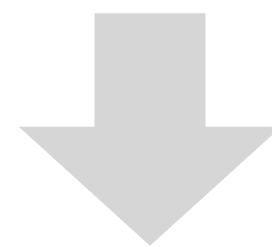
`TRUE & TRUE`



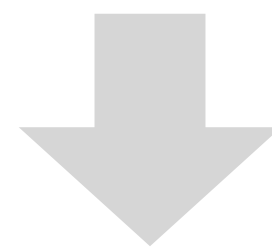
# Boolean operators

You can combine logical tests with `&`, `|`, `xor`, `!`, `any`, and `all`

`x > 2 & x < 9`



`TRUE & TRUE`



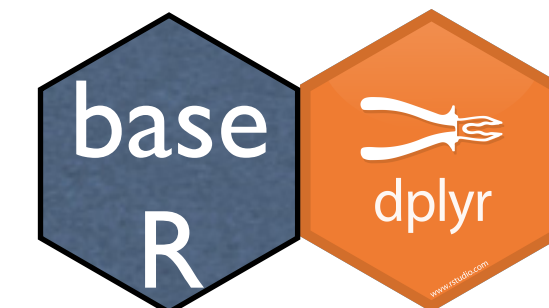
`TRUE`



# &

Are both condition 1 **and** condition 2 true?

expression	outcome
TRUE & TRUE	TRUE
TRUE & FALSE	FALSE
FALSE & TRUE	FALSE
FALSE & FALSE	FALSE



|

Is either condition 1 **or** condition 2 true?

expression	outcome
TRUE   TRUE	TRUE
TRUE   FALSE	TRUE
FALSE   TRUE	TRUE
FALSE   FALSE	FALSE



# XOR

Is either condition 1 **or** condition 2 true, **but not both**?

expression	outcome
xor(TRUE, TRUE)	FALSE
xor( TRUE, FALSE)	TRUE
xor( FALSE, TRUE)	TRUE
xor( FALSE, FALSE)	FALSE



!

## Negation

expression	outcome
!(TRUE)	FALSE
!(FALSE)	TRUE

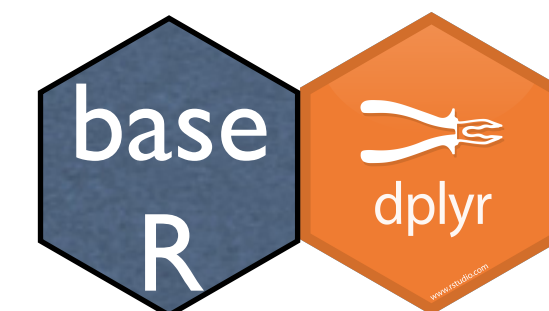
# Boolean operators



# Boolean operators

?base::Logic

<code>a &amp; b</code>	and
<code>a   b</code>	or
<code>xor(a, b)</code>	exactly or
<code>!a</code>	not



# Boolean operators

You can combine logical tests with `&`, `|`, `xor`, `!`, `any`, and `all`

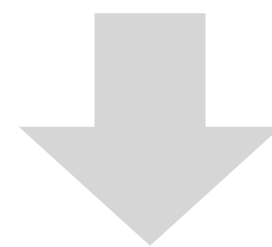
`x > 2 & x < 9`



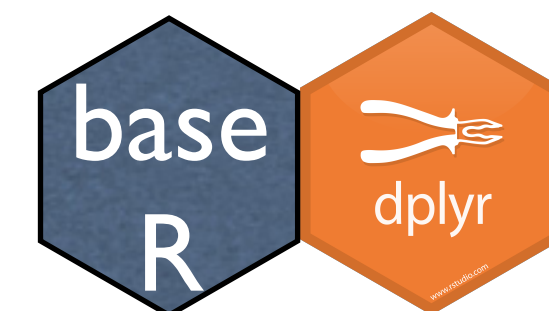
# Boolean operators

You can combine logical tests with `&`, `|`, `xor`, `!`, `any`, and `all`

`x > 2 & x < 9`



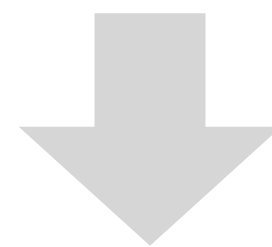
`TRUE &`



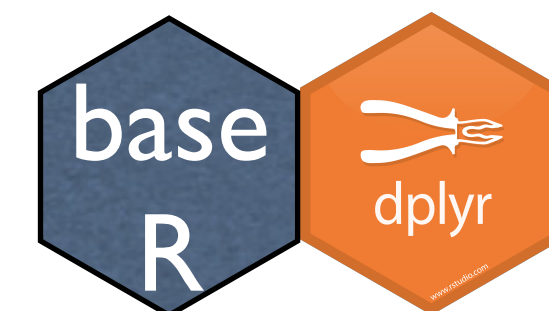
# Boolean operators

You can combine logical tests with `&`, `|`, `xor`, `!`, `any`, and `all`

`x > 2 & x < 9`



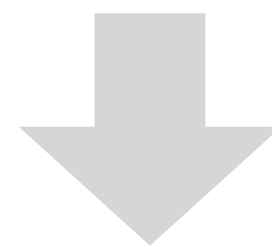
`TRUE & TRUE`



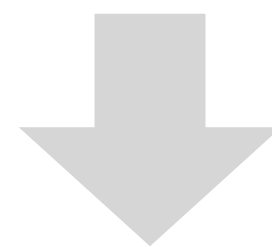
# Boolean operators

You can combine logical tests with `&`, `|`, `xor`, `!`, `any`, and `all`

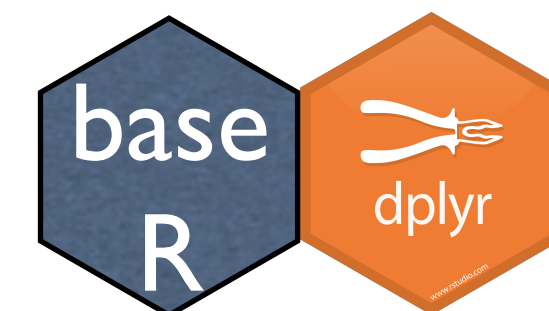
`x > 2 & x < 9`



`TRUE & TRUE`



`TRUE`



# &

Are both condition 1 **and** condition 2 true?

expression	outcome
TRUE & TRUE	TRUE
TRUE & FALSE	FALSE
FALSE & TRUE	FALSE
FALSE & FALSE	FALSE



|

Is either condition 1 **or** condition 2 true?

expression	outcome
TRUE   TRUE	TRUE
TRUE   FALSE	TRUE
FALSE   TRUE	TRUE
FALSE   FALSE	FALSE





# XOR

Is either condition 1 **or** condition 2 true, **but not both**?

expression	outcome
xor(TRUE, TRUE)	FALSE
xor( TRUE, FALSE)	TRUE
xor( FALSE, TRUE)	TRUE
xor( FALSE, FALSE)	FALSE

!

## Negation

expression	outcome
!(TRUE)	FALSE
!(FALSE)	TRUE

# filter()

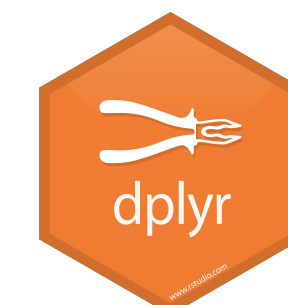
Extract rows that meet *every* logical criteria.

```
filter(beatles, birth==1940 & instrument == "guitar")
```

name	birth	instrument
John	1940	guitar
Paul	1942	base
George	1943	guitar
Ringo	1940	drums

→

John	1940	guitar
------	------	--------





# Your Turn

Modify the code below to filter out the rows for which birth is 1943 or instrument is drums

```
filter(beatles, birth==1940 & instrument == "guitar")
```

01:00




Base R

# Logical

You can subset with a logical vector of the same length as the dimension you are subsetting. Each element that corresponds to a TRUE will be returned.

```
beatles[c(FALSE, TRUE, TRUE, FALSE), ]
```

John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

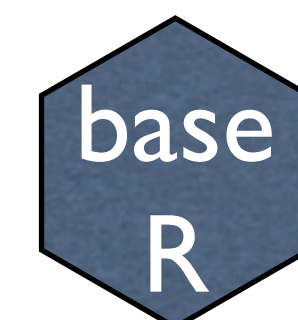


# Logical

You can subset with a logical vector of the same length as the dimension you are subsetting. Each element that corresponds to a TRUE will be returned.

```
beatles[c(FALSE,TRUE,TRUE,FALSE), ]
```

John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums





# Logical

You can provide a statement that **evaluates** to a logical to get something similar to a `dplyr filter()` statement.

```
beatles[beatles$birth == 1940, ]
```

```
beatles[c(TRUE, FALSE, FALSE, TRUE), ]
```

John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums





# Your Turn

Modify the code below to filter out rows where birth is 1943 or instrument is drums

```
beatles[beatles$birth == 1940, ]
```

01:00



More lists



# Quiz

What is the difference between an atomic vector and a list?



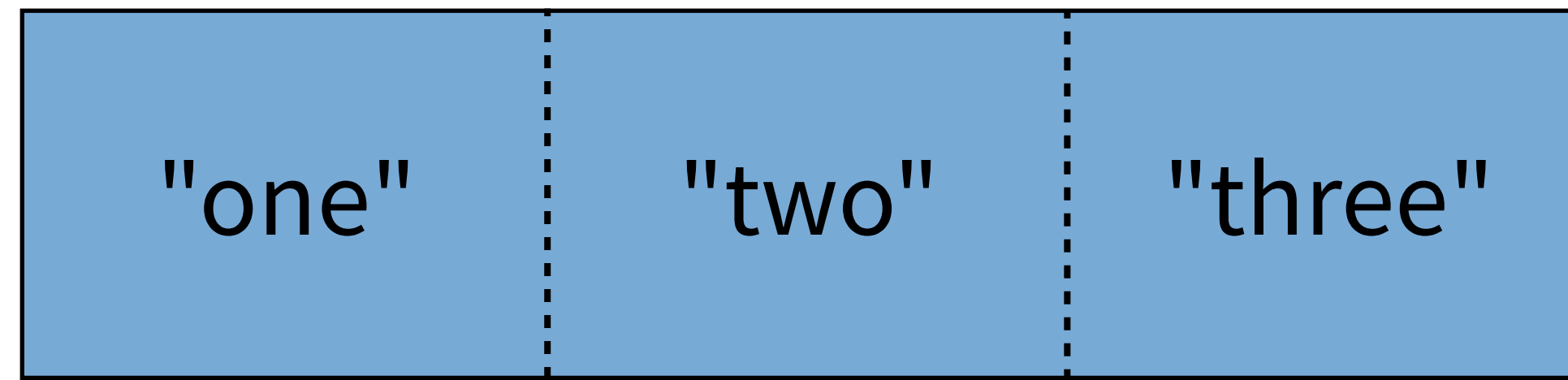
Atomic Vector



type



Atomic Vector

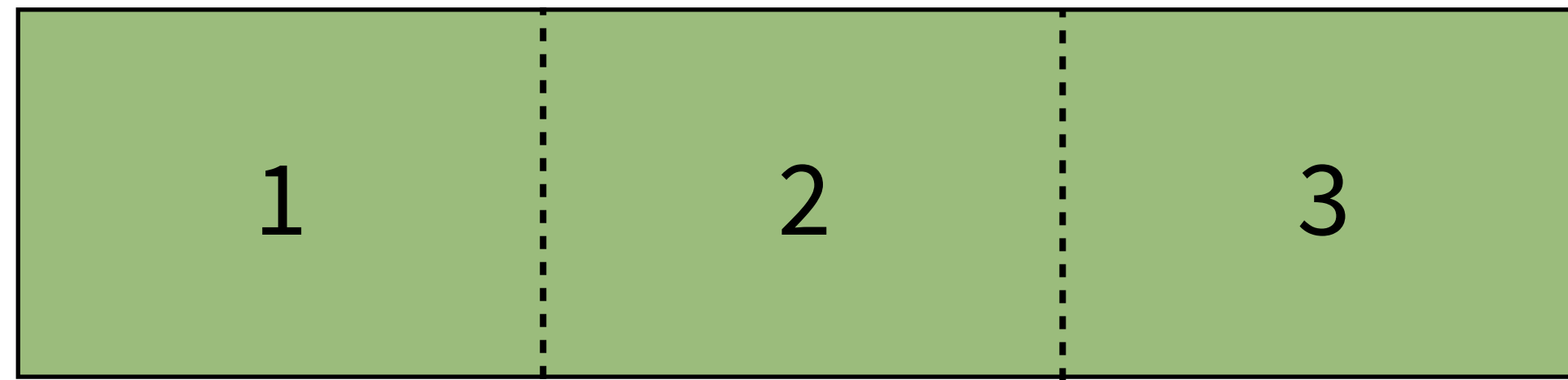


character



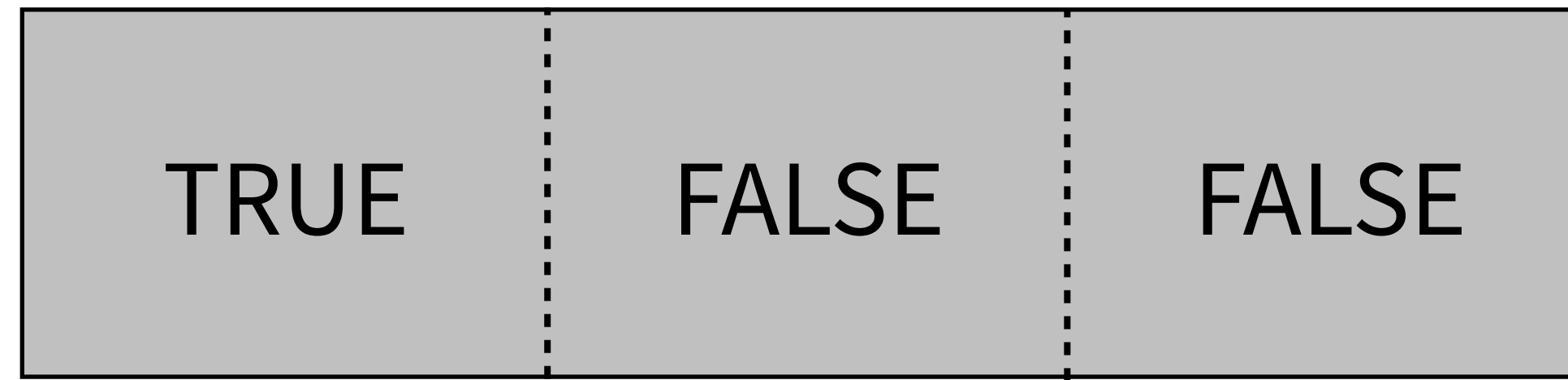


Atomic Vector



double

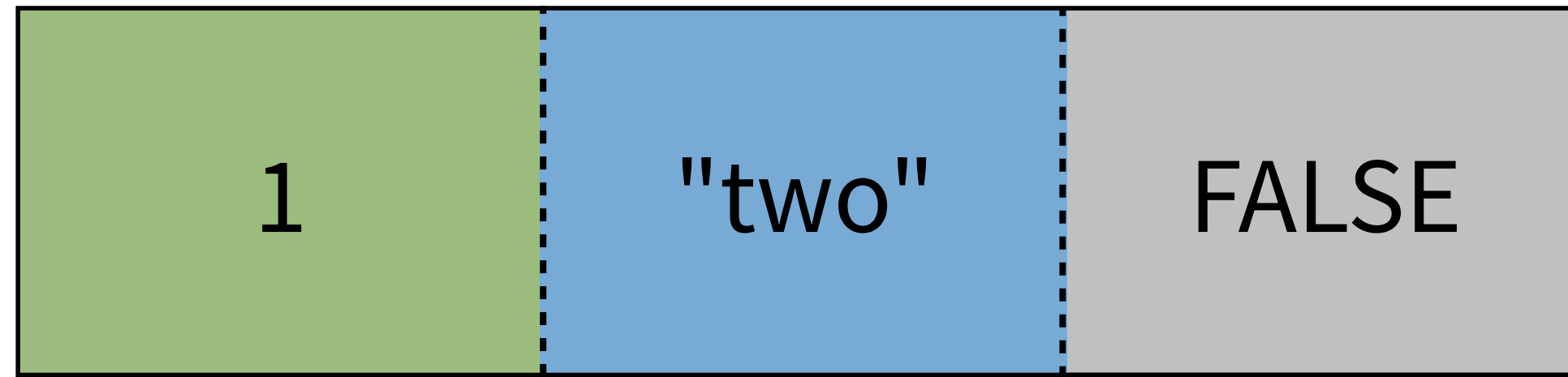
Atomic Vector



logical

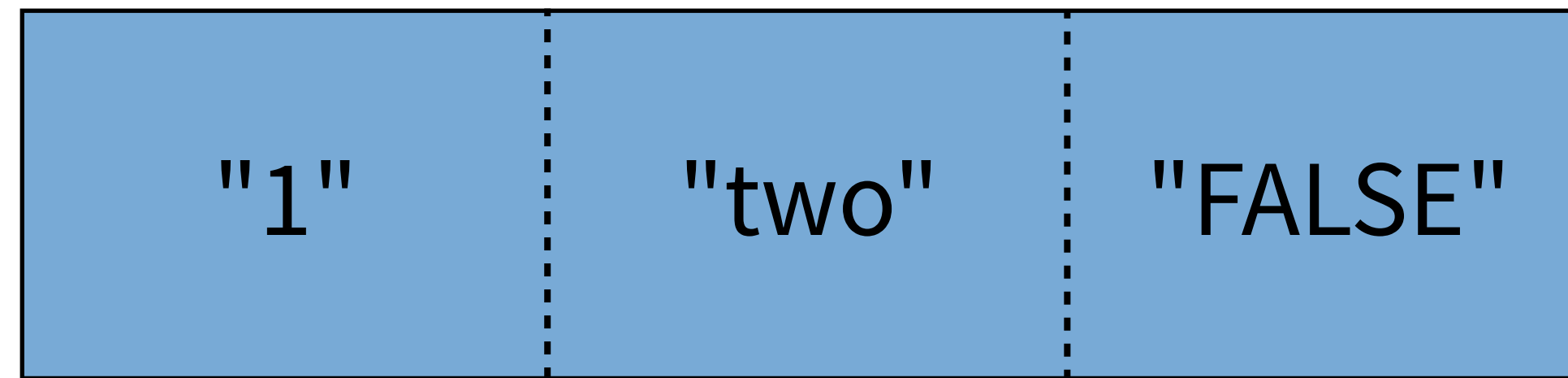


Atomic Vector



?

Atomic Vector



character

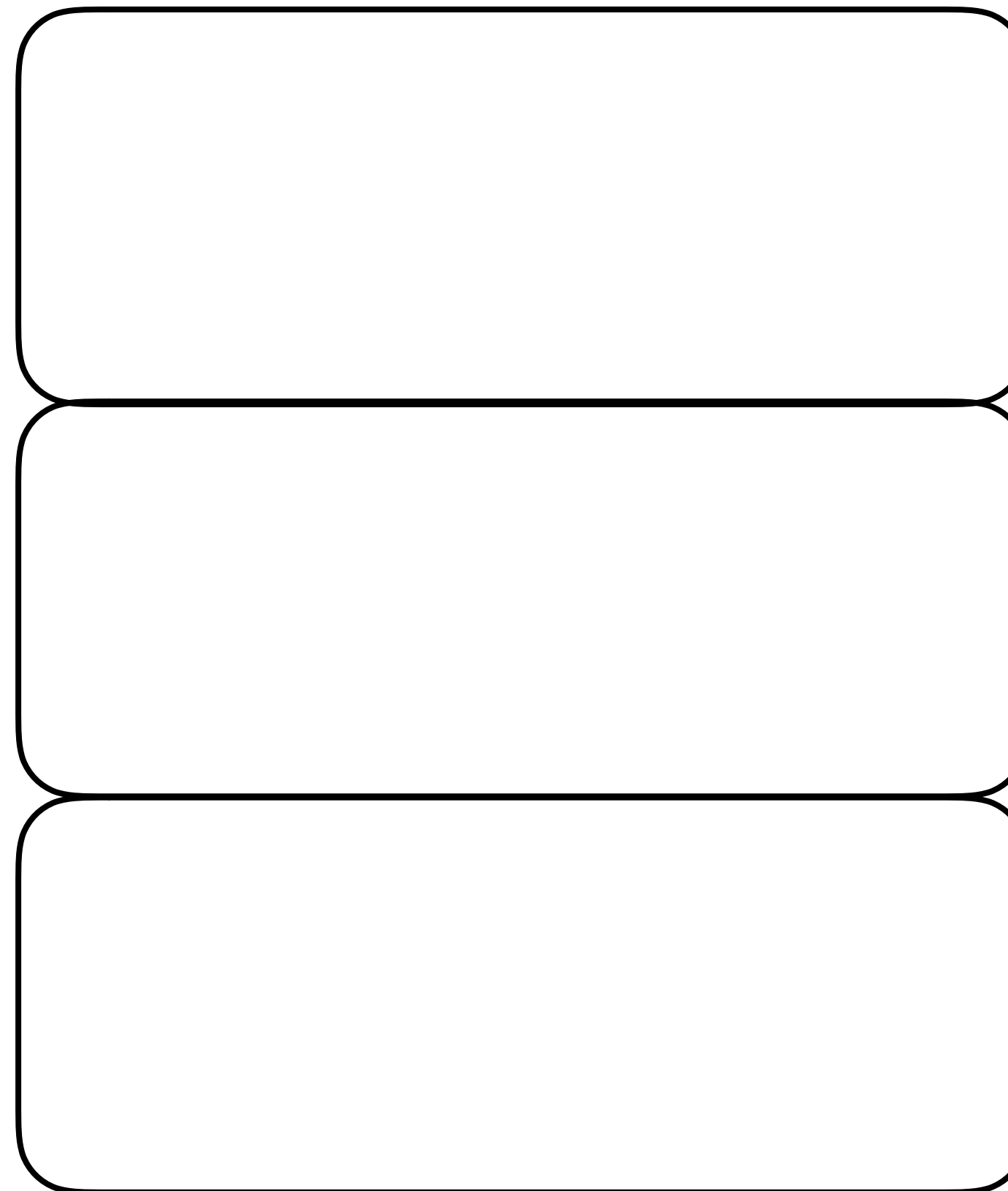


Atomic Vector



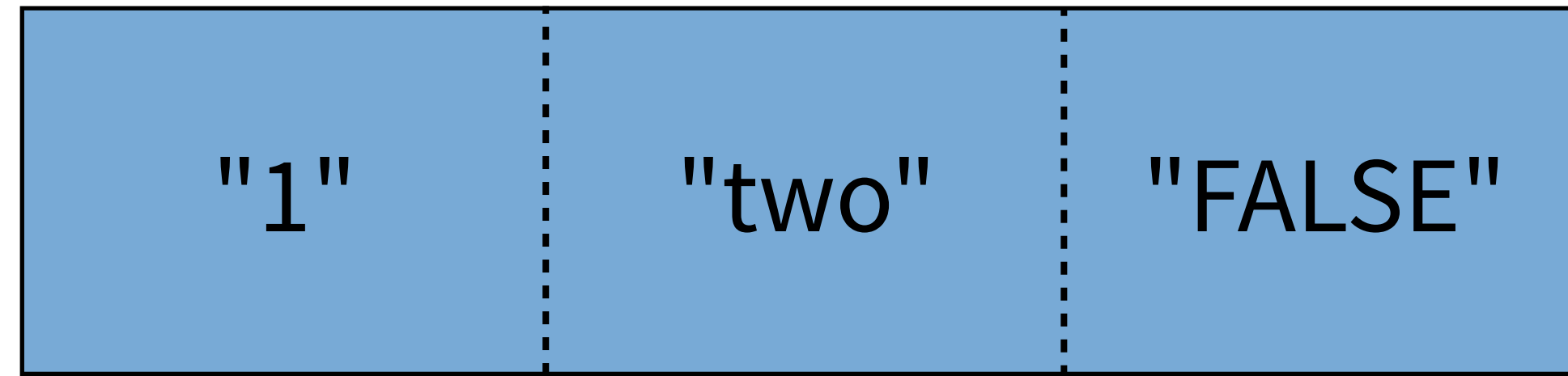
type

List



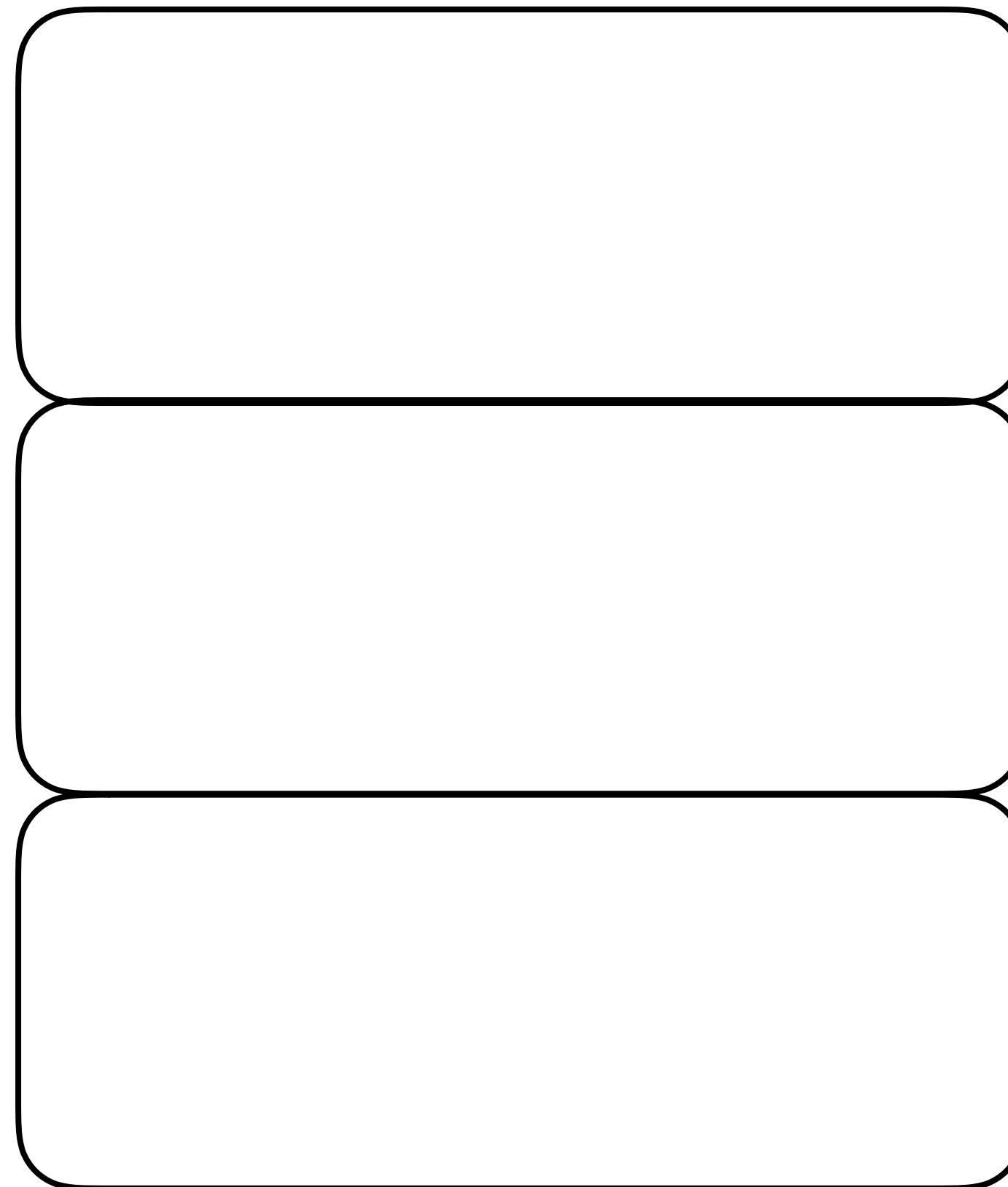


Atomic Vector

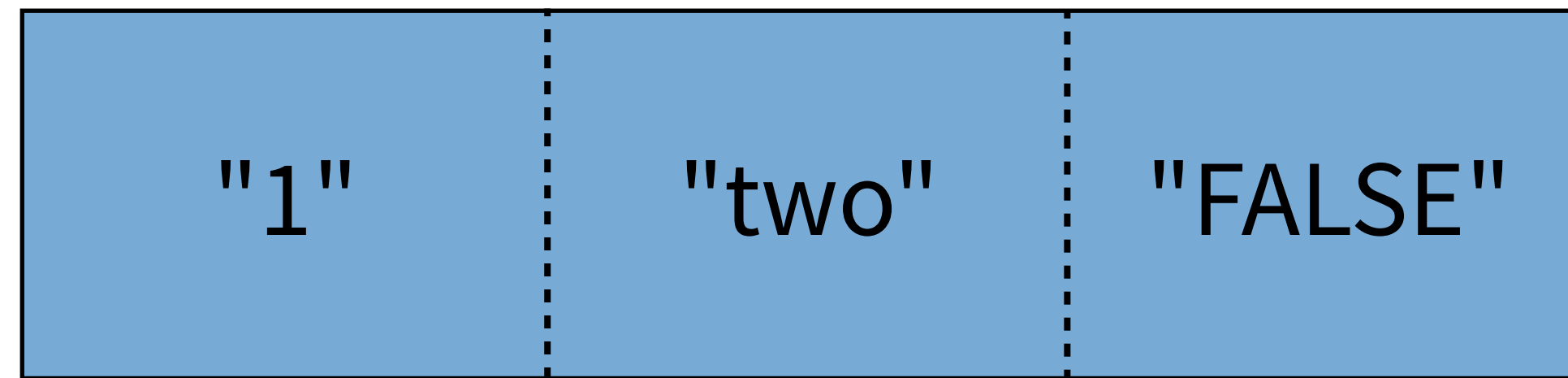


character

List

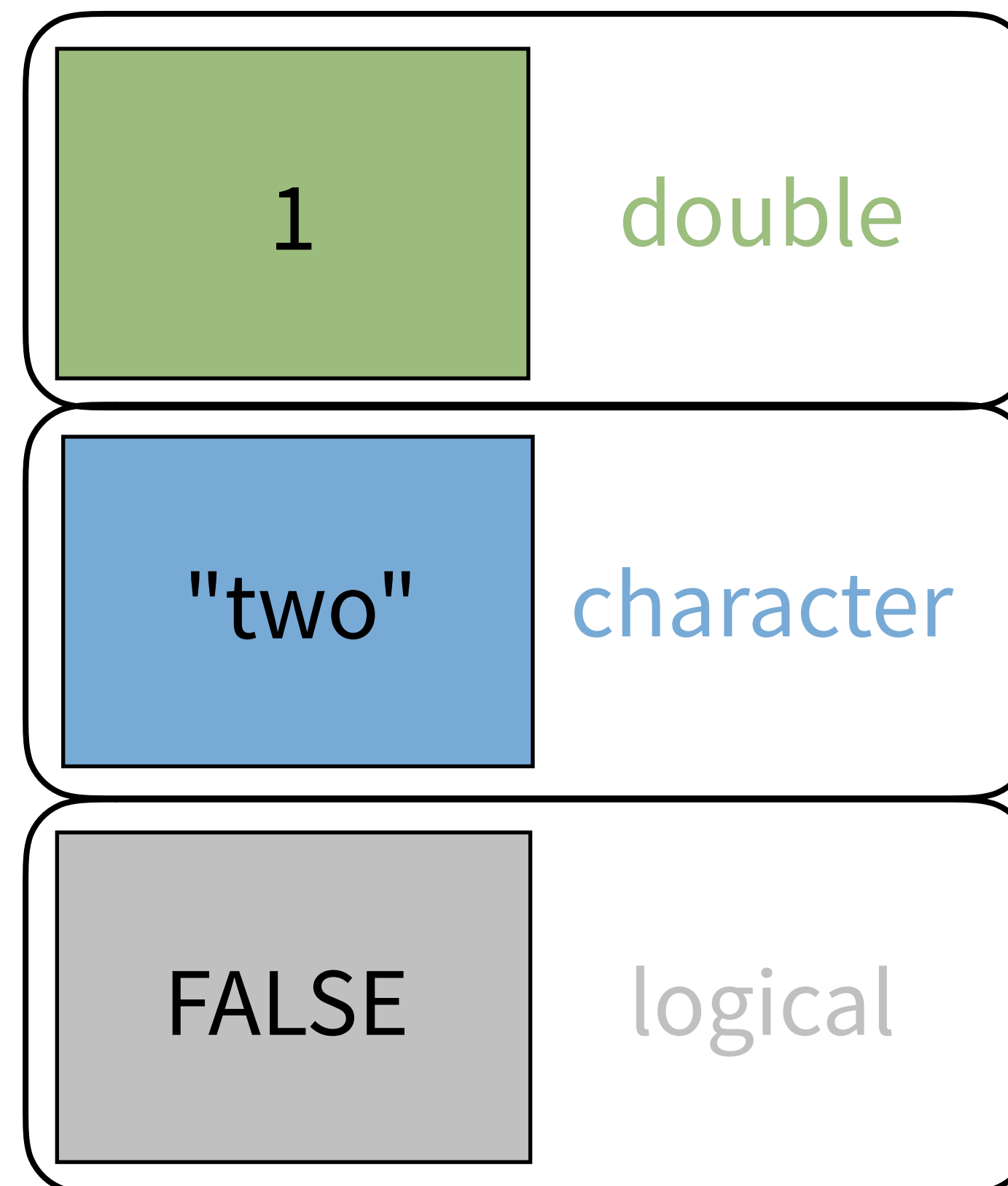


Atomic Vector

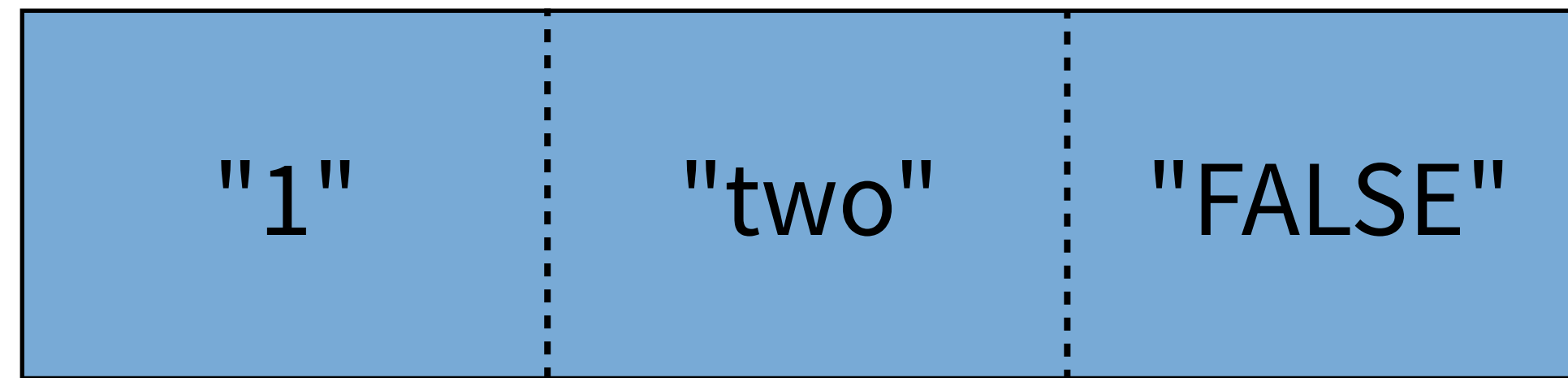


character

List

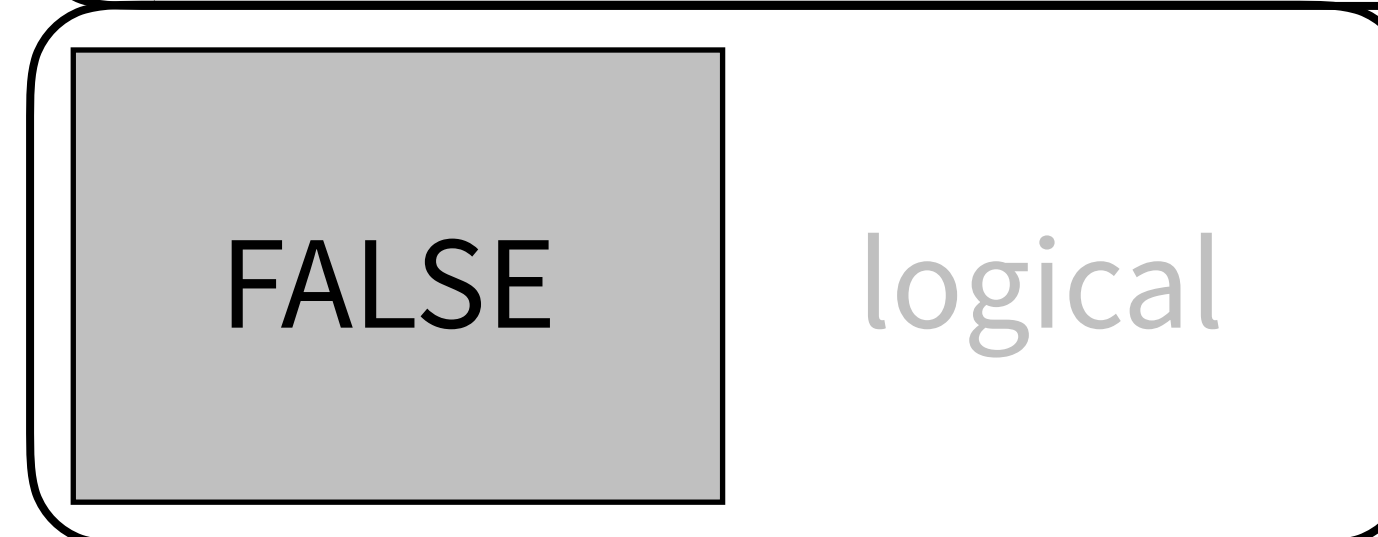
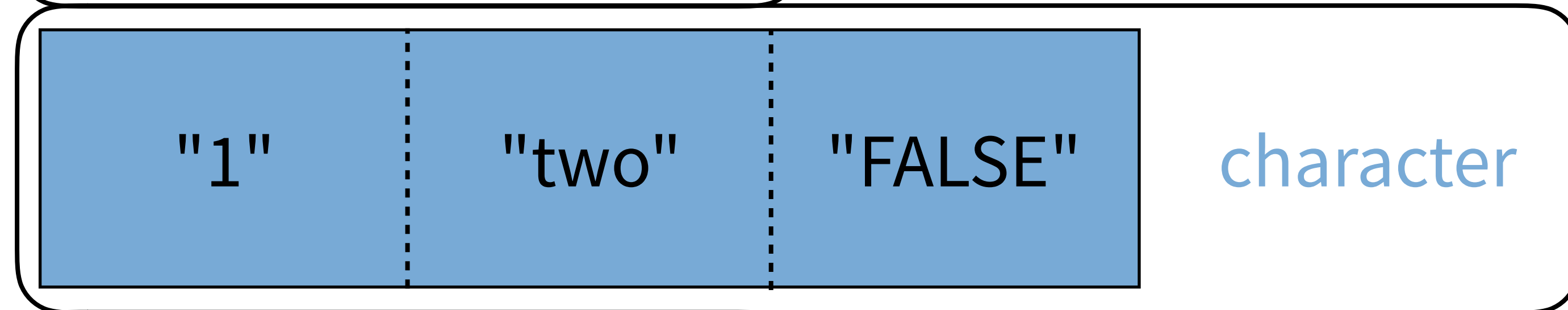
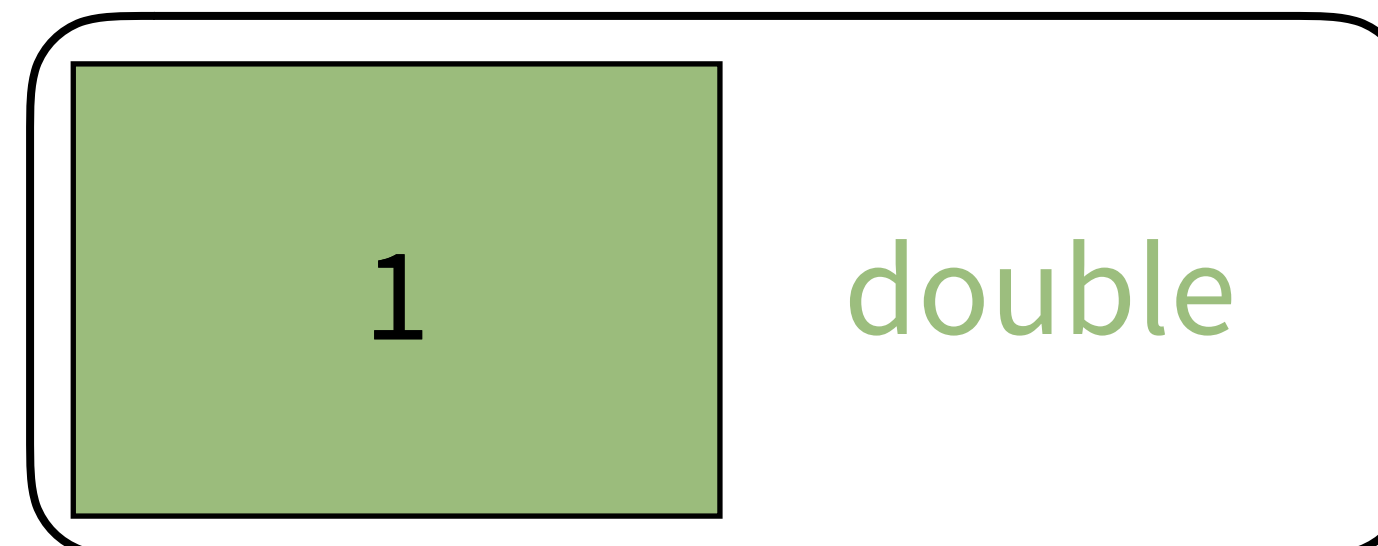


Atomic Vector

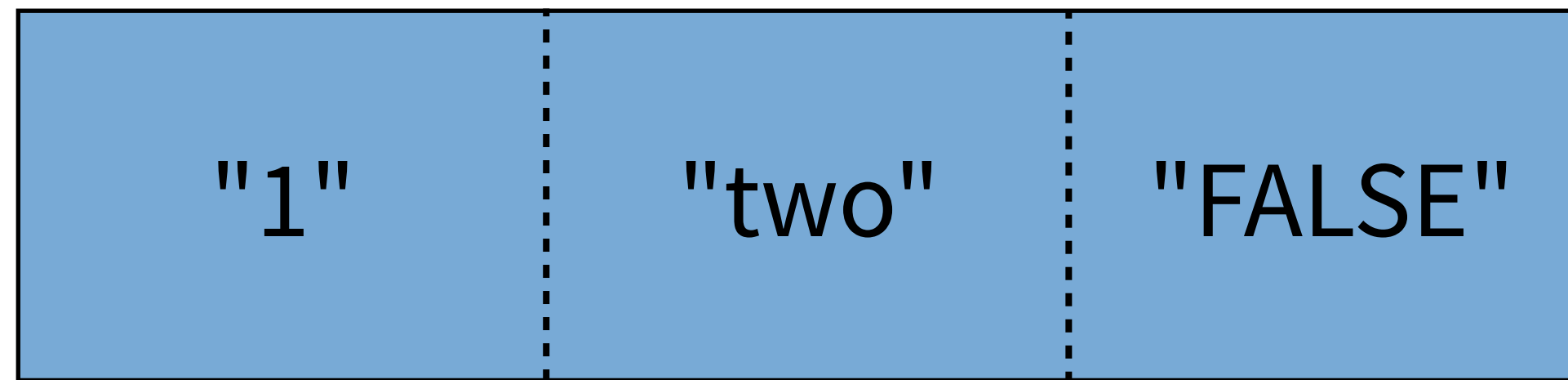


character

List

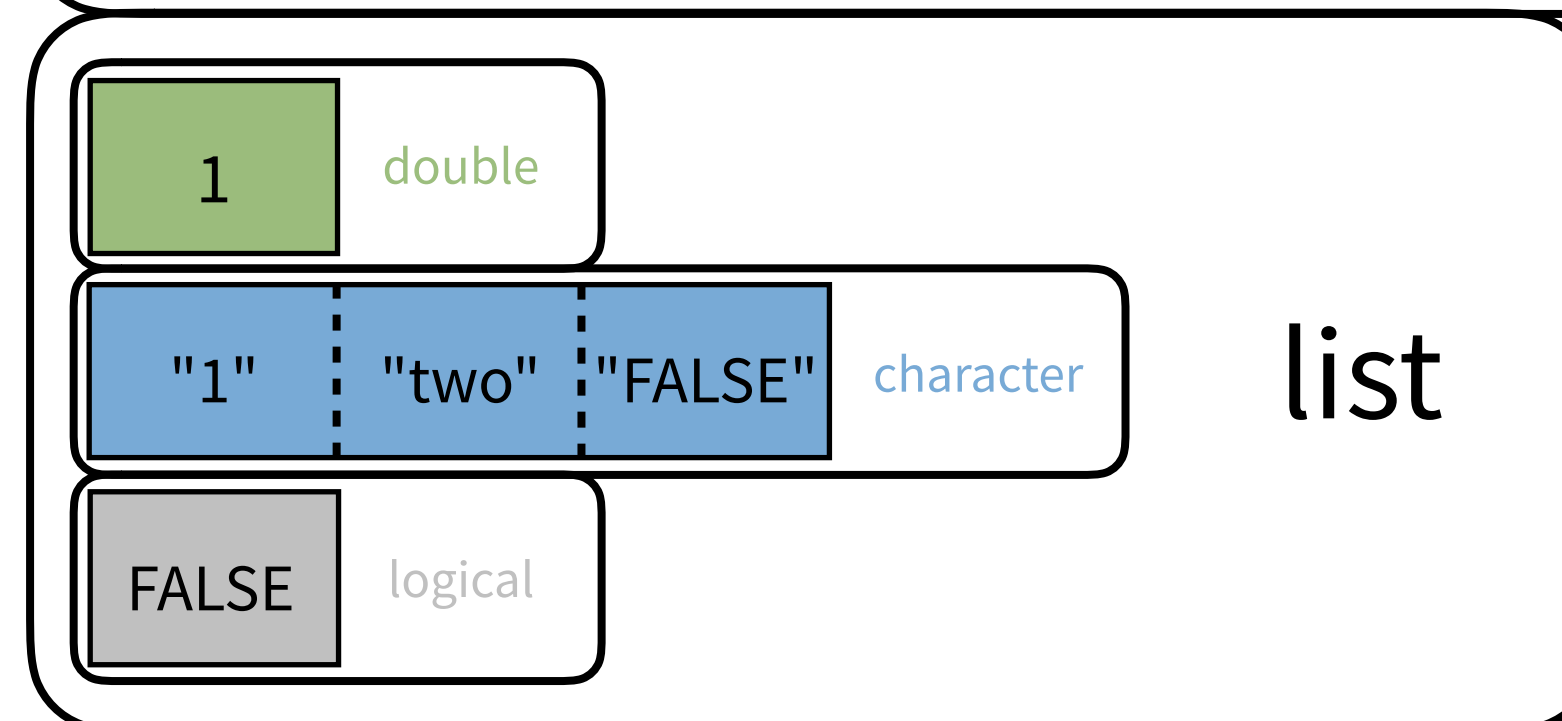
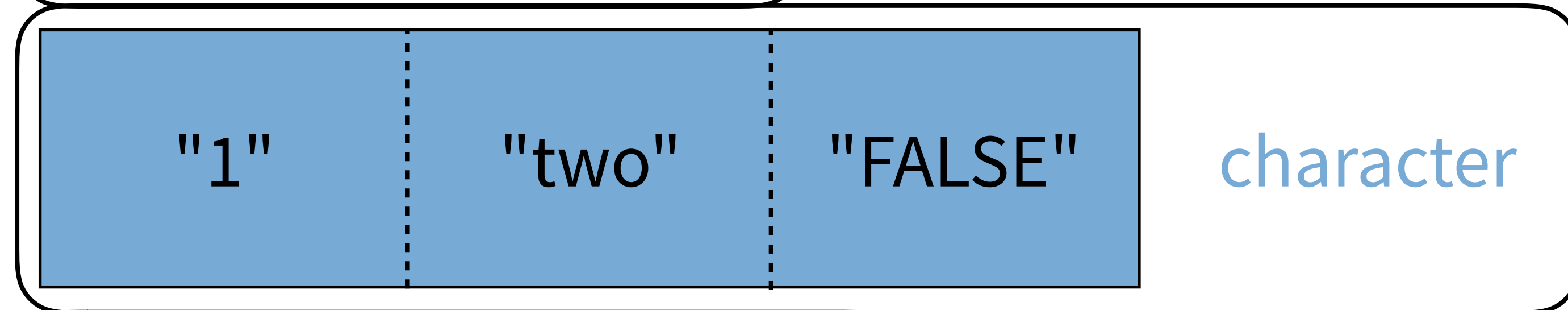
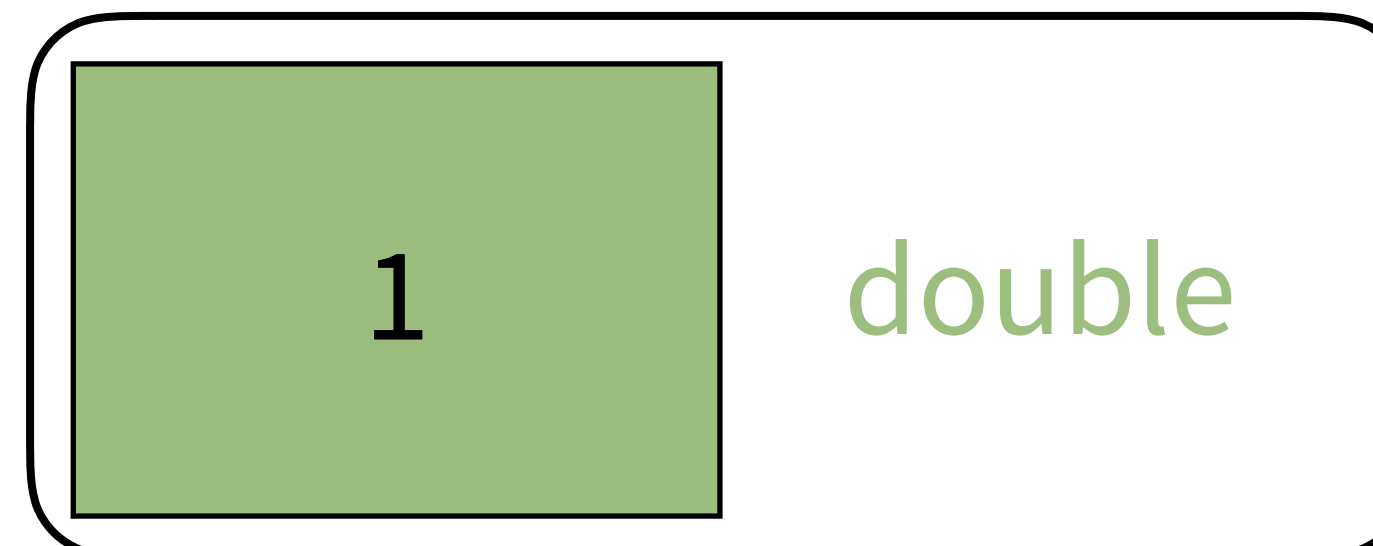


Atomic Vector



character

List





# Your Turn 1

Here is a list:

```
a_list <- list(num = c(8, 9),  
              log = TRUE,  
              cha = c("a", "b", "c"))
```

Here are two subsetting commands. Do they return the same values? Run the code chunks to confirm

```
a_list["num"]
```

```
a_list[["num"]]
```



```
a_list["num"]
```

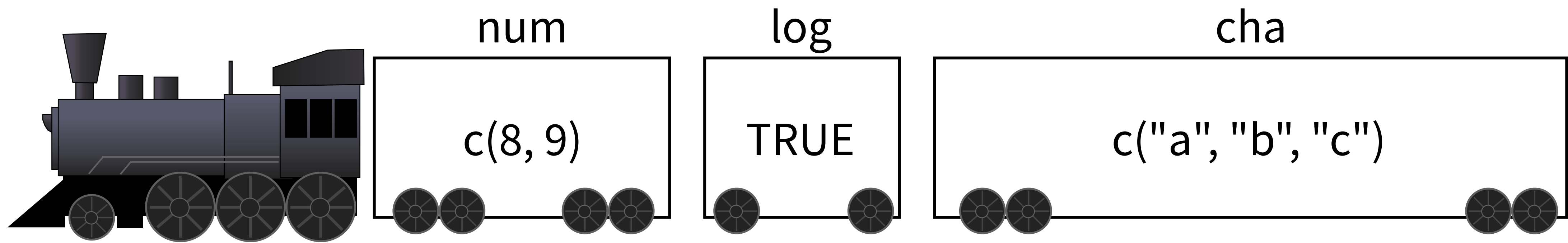
```
$num  
[1] 8 9
```

**A list**  
(with one element named  
num that contains an  
atomic vector)

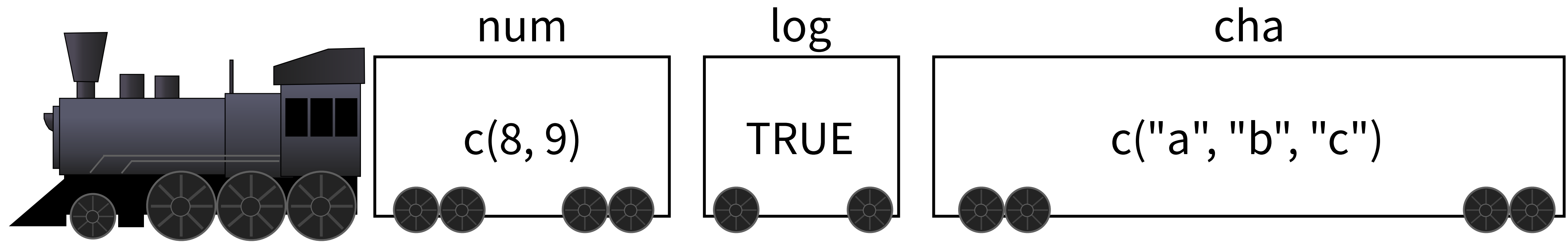
```
a_list[["num"]]
```

```
[1] 8 9
```

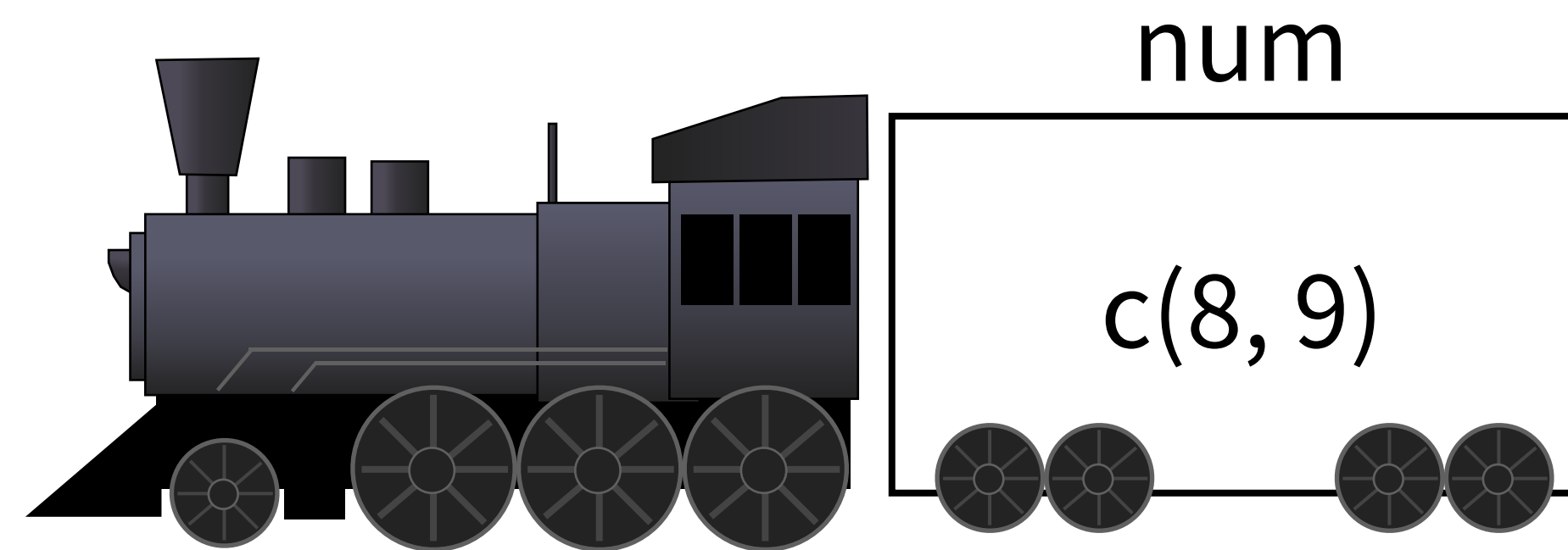
**An atomic vector**

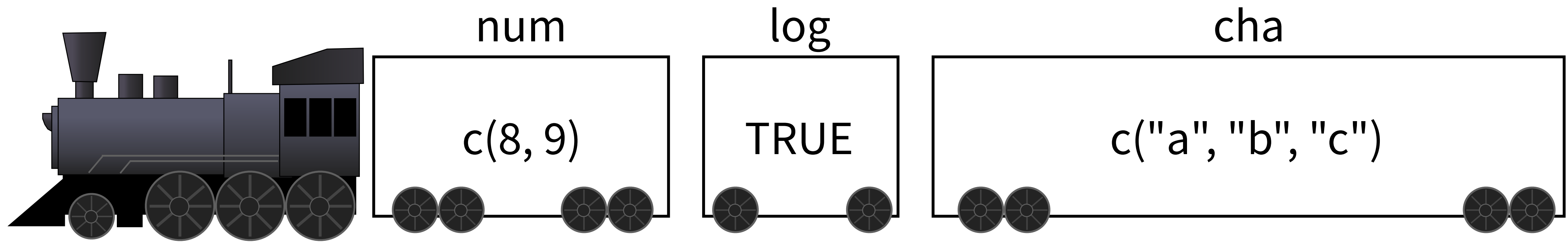


```
a_list <- list(num = c(8,9), log = TRUE, cha = c("a", "b", "c"))
```

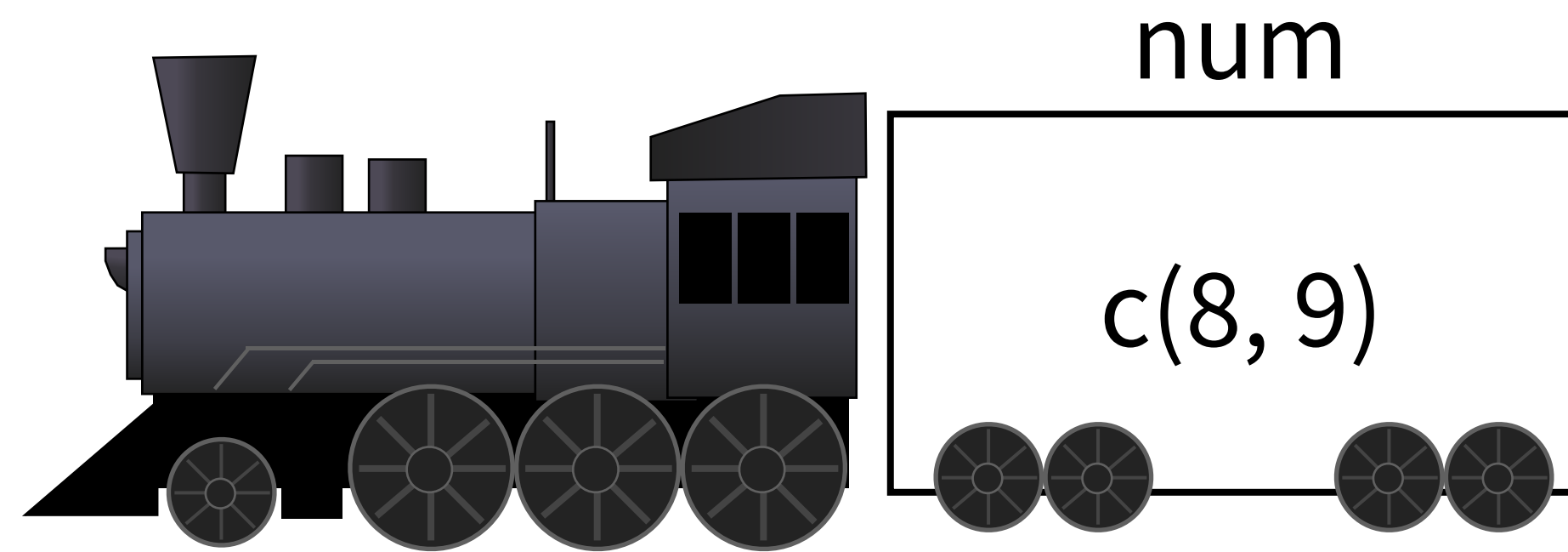


```
a_list["num"]
```



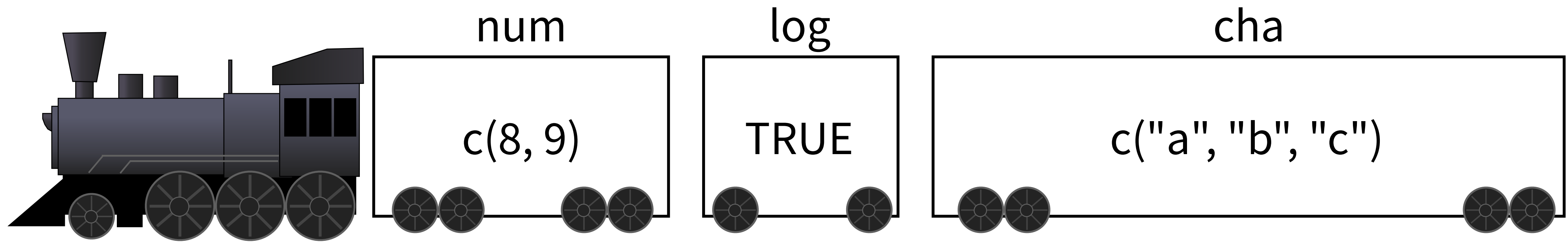


```
a_list["num"]
```

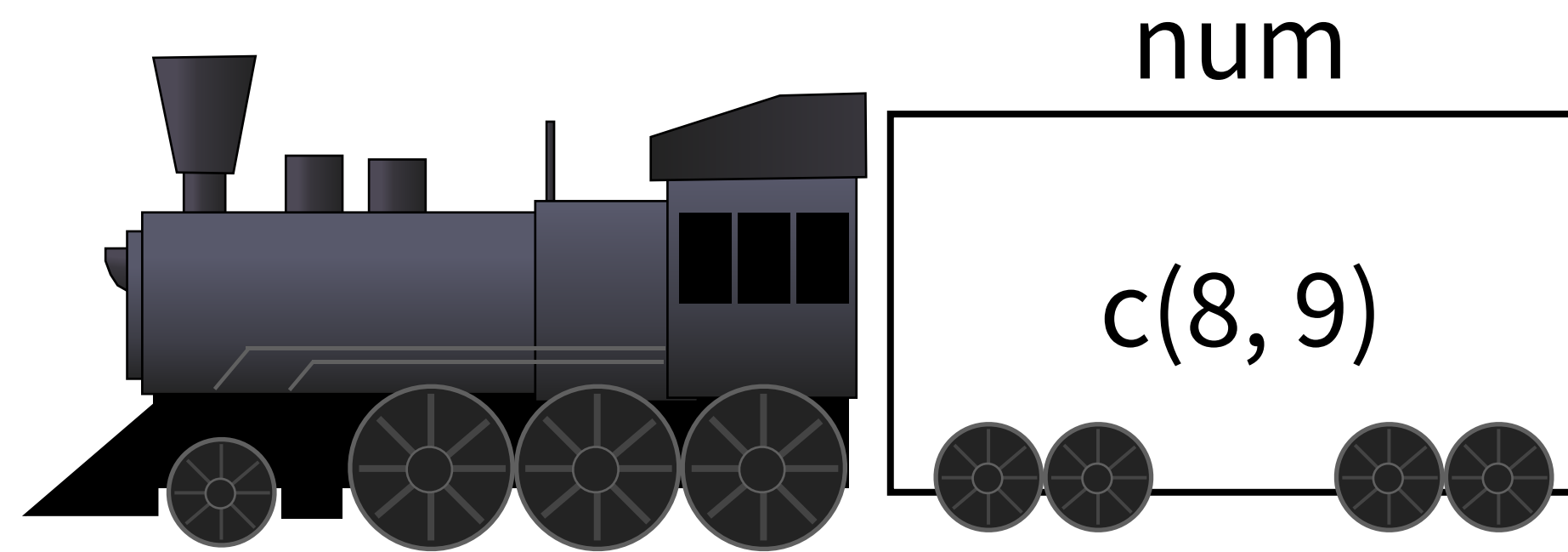


```
a_list[["num"]]
```

`c(8, 9)`



```
a_list["num"]
```



```
a_list[["num"]]
```

`c(8, 9)`

```
a_list$num
```

`c(8, 9)`