

R Lesson 2: Objects and Data Structures

vanderbi.lt/r

Steve Baskauf



Preliminaries



Common types of data

- **character**, e.g. "Fred" or "!@#ts23" (in quotes)
- **numeric**, e.g. 15 or 6.02 (no quotes)
- **logical**, TRUE or FALSE (all caps, no quotes)

Object name recommendations

- Be descriptive (what the object is or does)
- **snake_case** (underscores) is commonly used:
 - `ordinary_relational_processes`
- camelCase is sometimes used:
 - `bookList, alphabetizeParticipants`
- We can use the term **variable** to refer to named objects
- R doesn't know what a name "means". A meaningful name helps human readers of the code.

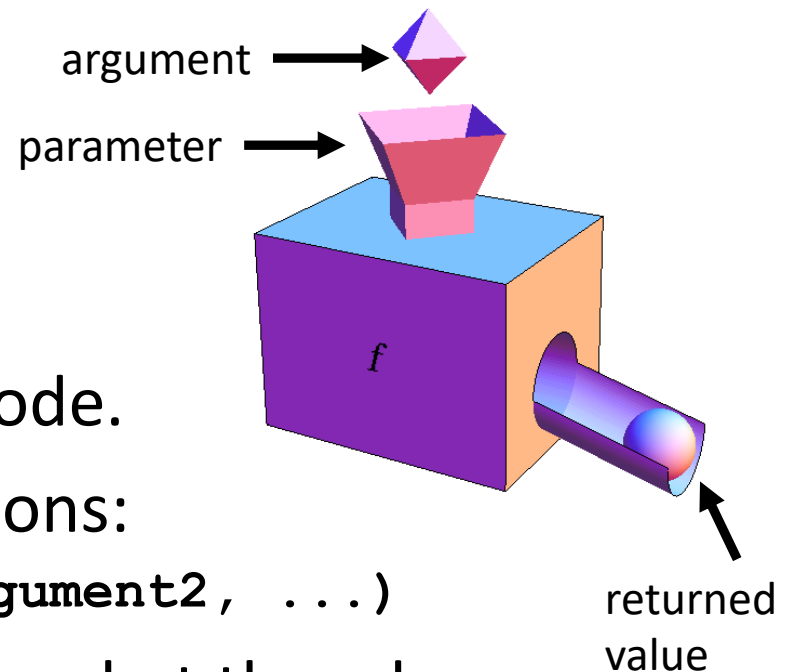
Assigning a value to an object

- You can **assign** a value to an object using **<-** (similar to a left arrow)
- Examples:
 - `name <- "Steve"` (creating a character object)
 - `my_number <- 6.02` (creating a numeric object)
- Using the equals sign (=) is allowed, but not recommended.
- alt-minus is an RStudio shortcut to generate <-

"Printing" the value of an object

- R does not have a "print" command.
- entering the name of an object (or expression) in the console evaluates and displays its value

Functions



- A **function** defines a block of code.
- We pass **arguments** into functions:
`function_name(argument1, argument2, ...)`
- Functions are usually named by what they do.
Example:

```
my_latte <- make_latte(beans, milk, water)
```

- Functions can be:
 - built-in to R
 - defined by you in your code
 - defined by somebody else in a package

Using a function

- We don't have to know anything about the code that makes a function work. We just need to know:
 - What the function does
 - What arguments to put into it
 - What the function will output
- Examples:
 - `sqrt(2)` (evaluate and display)
 - `x <- sqrt(3)` (evaluate and assign to an object)

Vectors



Vectors are king in R

vector named `animal`

"frog"	"spider"	"worm"	"bee"
<code>animal[1]</code>	<code>animal[2]</code>	<code>animal[3]</code>	<code>animal[4]</code>

- A **vector** is the most common kind of data structure in R.
- Vectors contain a sequence of the **same type** of data.

Creating vectors

- We commonly use the **construct** function to make vectors:

```
number_vector <- c(1, 3, 6, 10, 15)
```

```
animal <- c("frog", "spider", "worm", "bee")
```

- We can also generate a **sequence** of numbers:

```
number_range <- 3:9
```

```
count_down <- 10:0
```

```
go_negative <- 5:-3
```

- The generated sequence is just another vector!
- (Python users: note the range includes the final value)

Knowing what's going on with a vector

- display it in console
- examine its value in the environment data pane
- examine its properties:
 - `length(animals)` (how many items)
 - `mode(animals)` (type of data in vector)

Referencing parts of vectors

- Referencing a **single item**:
`animal[3]` (displays the third item)
`animal[2] <- "arachnid"` (assigns "arachnid" to the 2nd item)
- Referencing a **range of items** (subvector):
`animal[2:4]` (the range 2:4 is actually a vector itself)
- (Python users: R vectors are "1 based"; the first item is numbered 1, not 0. Also, the range includes the final value.)

Single item objects are vectors, too.

- Surprisingly, a single data item assigned to an object is also a vector. We can see this if we ask its length as if it were a vector:

```
an_item <- "some character string"  
length(an_item)
```

- We can reference the single item using vector notation:

```
an_item[1]
```

Operations on vectors

- Many functions work equally well for a single item or a multi-item vector (since they are both vectors):
`number_vector <- c(1, 3, 6, 10, 15)`
`sqrt(number_vector)`
- When operations are performed on vectors, they generally are performed on **all items** in the vectors.

```
> a <- c(10, 30, 100)
> b <- c(5, 10, 20)
> c <- a/b
> c
[1] 2 3 5
```

More complicated things are also vectors

- A **matrix** is a vector that has been assigned two dimensions
- An **array** is a vector that has been assigned any number of dimensions
- As forms of vectors, matrices and arrays can only consist of one kind of data.
- Example:

```
a_vector <- c(1.1, 1.2, 2.1, 2.2, 3.1, 3.2)
```

```
a_matrix <- matrix(a_vector, 2, 3)
```


Missing data indicators

- R's built-in indicators for **missing data**:
 - **NA** ("not available") means there is a value, but it's missing; length =1
 - NULL means no value; length=0

```
vector_with_missing <- c(1, 2, NA, 3)
```

- NA will prevent some calculations. Example:

```
mean(vector_with_missing)
```

- NA can be used for missing data in tables instead of blank cells

Other important data structures



Lists

list named thing

name	<i>fruitKind</i>	<i>euler</i>	<i>vectorData</i>	<i>curse</i>
value	"apple"	2.71828	animal	"!@#\$\$%"

reference value by position `thing[[1]]` `thing[[2]]` `thing[[3]]` `thing[[4]]`

reference value by name `thing$fruitKind` `thing$euler` `thing$vectorData` `thing$curse`

- Like vectors, **lists** are one-dimensional data structures.
- However, lists can be **heterogeneous** (contain more than one kind of data object)
- It is typical to give names to values of a list.

Creating a list

- Lists are created using the `list()` function:

```
thing <- list(fruit_kind="apple",  
             euler=2.71828,  
             vector_data=animal,  
             curse="!@#$%")
```

- This list contains character strings, a number, and a vector.
- Values can be assigned names as they are added to the list

Viewing contents of a list

The screenshot shows the RStudio interface. On the left, the Environment pane displays a table of objects:

Name	Type	Value
thing	list [4]	List of length 4
fruitKind	character [1]	'apple'
euler	double [1]	2.71828
vectorData	character [4]	'frog' 'spider' 'worm' 'bee'
curse	character [1]	'!@#\$\$'

On the right, the Environment pane shows the details for the 'thing' object:

Data	Value
thing	List of 4
Values	chr [1:4] "frog" "spider" "worm" "bee"

Red arrows point from the 'thing' row in the left pane to the 'thing' row in the right pane, and from the 'Values' row in the right pane back to the 'thing' row in the left pane. Red text annotations are present:

- the list shows up in the workspace summary** (pointing to the 'thing' row in the right pane)
- clicking on it brings up details in the upper left pane** (pointing to the 'thing' row in the left pane)

- You can see what's in a list by clicking on its name in the workspace summary in the Environment pane

Referencing list items

- List items can be referenced by:
 - **position** using double square brackets and the index number

```
thing[[2]]
```

- **name** using a dollar sign and the name string

```
thing$curse
```

Clearing the contents of a pane

- Click on the little broom near the top of the pane
- The view in the pane will be cleared
- In the case of the Environment pane, the values will also be cleared.

Data frames

data frame named `organismInfo`

column name

group	animal	numberLegs
"reptile"	"frog"	4
"arachnid"	"spider"	8
"annelid"	"worm"	0
"insect"	"bee"	6

`organismInfo[2,1]`

vector

`organismInfo[4,3]`

`organismInfo$animal[4]`

- **Data frames** are essentially tables
- The **column values** are like **vectors**
- The **set of columns** is like a **list**

Making a data frame from vectors

- First make the named vectors

```
group <- c("reptile", "arachnid", "annelid",  
"insect") # vector of strings
```

```
animal <- c("frog", "spider", "worm", "bee")
```

```
number_legs <- c(4, 8, 0, 6) # vector of numbers
```

- Then put the vectors into a data frame

```
organism_info <- data.frame(group, animal,  
number_legs)
```

- The vector names will be used for the column names

Viewing contents of a data frame

The screenshot shows the RStudio interface with the 'organismInfo' data frame displayed as a table in the Environment pane. The table has three columns: 'group', 'animal', and 'numberLegs'. The data is as follows:

	group	animal	numberLegs
1	reptile	frog	4
2	arachnid	spider	8
3	annelid	worm	0
4	insect	bee	6

The Environment pane shows the data frame 'organismInfo' with 4 observations and 3 variables. The values for each variable are:

Variable	Values
animal	chr [1:4] "frog" "spider" "worm" "bee"
group	chr [1:4] "reptile" "arachnid" "annelid" "insect"
numberLegs	num [1:4] 4 8 0 6

click on the data frame name here to see it displayed as a table here

- Click on the name of the data frame in the Environment pane
- The contents will be displayed as a table

Referring to parts of a data frame

- Since the columns are like list items, we can refer to them by name:

organism_info\$animal

- Individual cells can be referenced by:

- row and column

organism_info[2,1]

- column name and position in column

organism_info\$animal[4]

Examining a data frame

- **head()** shows the first 6 rows
- **tail()** shows the last 6 rows
- **names()** returns the column names
- **str()** describes the structure of the data frame with information about each column

Loading data from files



Tabular data in delimited files

- **Delimited files** are text files where values are separated by some text character and lines are separated by **newline** characters (i.e. "hard returns").
- Most common type of delimited file: **CSV** (comma separated values)
- Also used: TSV (tab separated values)
- Delimited files are much simpler than Excel files and are commonly used for archiving data.
- CSV files can be made by exporting from Excel

Reading delimited files into data frames

- There are several ways to read data from CSV files into R:

- by a **file path** (platform-dependent)

```
my_data_frame <- read.csv("~/test.csv") (Mac)
```

```
my_data_frame <- read.csv("c:\\temp\\test.csv") (Windows)
```

- by a **file-choosing dialog**

```
my_data_frame <- read.csv(file.choose())
```

- by a **URL**

```
my_data_frame <-  
read.csv("https://gist.githubusercontent.com/baskaufs/1a7  
a995c1b25d6e88b45/raw/4bb17ccc5c1e62c27627833a4f25380f27d  
30b35/t-test.csv")
```

Controlling the import process

- You can specify if the file has a **header row** (labels) using the **header** key (default value is TRUE)
- You can specify the **separator** if it's different from comma using the **sep** key (default value is comma)
- `\t` is the escaped value for a tab character
- Example:

```
nls_ds1 <- read.csv(file.choose(),  
                    header = TRUE,  
                    sep = "\t")
```


Homework: Nashville schools data

1. What does R do when column headers have spaces in them?
2. Display the values in the zip code column
3. How many values are there in the zip code column?
4. Calculate the number of students in each school by adding the values in the male and female columns
5. Calculate the fraction of students that are white in each school
6. Calculate the average fraction of white students by school