

Welcome & Types in R

Lecture 01

Dr. Colin Rundel

Course Details

Course Team

Instructor

- Dr. Colin Rundel
 - colin.rundel@duke.edu / cr173@duke.edu / rundel@gmail.com
 - Office hours: 204 Old Chem or Zoom - Time TBD or by appointment

TAs

- Lynn Kremers
- Grady Purcell

Course website(s)

- GitHub pages - <https://sta323-sp26.github.io>
 - HTML, PDF, and qmds of Slides
 - Readings and other notes
- Canvas - <https://canvas.duke.edu/courses/73157>
 - Announcements
 - Gradebook

Labs

- Attendance is expected - must attend lab you are enrolled in
- Opportunity to work on course assignments with TA support
- Labs will begin next week - Monday (1/12)

Assessment

This course is graded 100% on your coursework (there are no exams).

We will be assessing you based on the following:

Assignment	Type	Value	n	Assigned
Homeworks	Team	30%	5/6	~ Every other week
Midterms	Individual	40%	2	~ Week 6 and 14
Project	Team	10%	1	~ Week 10
Quizzes	Individual	20%	~10-12	~ Wee

Teams

- Roughly biweekly homework assignments
- Open ended, ~5 - 15 hours of work
- Peer evaluation after completion
- Expectations and roles:
 - Everyone is expected to contribute equal *effort*
 - Everyone is expected to understand *all* code turned in
 - Individual contribution evaluated by multiple metrics (peer evaluation, commits, etc.)

Collaboration policy

- Only work that is clearly assigned as team work should be completed collaboratively (Homeworks + Project).
- Individual assignments (Midterms) must be completed individually, you may not directly share or discuss answers / code with anyone other than the myself and the TAs.
- On Homeworks you should not directly share answers / code with other teams, however you are welcome to discuss the problems in general and ask for advice.

Sharing / reusing code / AI policy

- We are aware that a huge volume of code is available on the web, and many tasks may have solutions posted.
- Unless explicitly stated otherwise, this course's policy is that you may make use of any online resources (e.g. Google, StackOverflow, etc.) but you must explicitly cite where you obtained any code you directly use or use as inspiration in your solution(s).
- Any recycled/copied code that is not explicitly cited will be treated as plagiarism, regardless of source.
- The same applies to the use of LLM like ChatGPT, Claude, Gemini, or GitHub Copilot - you are welcome to make use of these tools as the basis for your solutions but you must cite the tool when using it for *significant* code generation.

Brief thoughts on AI tools

- AI tools are not a replacement for understanding the material, but they can be a tool to help you understand the material.
- Reading code and writing code are skills that take time and practice to develop - both are essential.
- Nature of the tools is changing rapidly - Autocomplete vs ChatBots vs Agentic
 - Google Gemini and GitHub Copilot have “free” student options

Academic integrity

To uphold the Duke Community Standard:

- I will not lie, cheat, or steal in my academic endeavors;
- I will conduct myself honorably in all my endeavors; and
- I will act if the Standard is compromised.

Course Tools

Accessing RStudio Workbench

To reduce friction, the preferred method is to use the department's RStudio server(s).

To access RStudio/Posit Workbench:

1. Navigate to <https://rstudio.stat.duke.edu>
2. Log-in with your Duke NetID and password.

If off campus, use the VPN to create a secure connection from your computer to Duke. If you are on campus, be sure you

DSS RStudio alternatives

If you cannot access RStudio via the DSS servers:

- Make sure you are on authenticated Duke network (e.g. DukeBlue or VPN)
- Make sure you are not using a custom DNS server
 - e.g. [1.1.1.1](#) or [8.8.8.8](#)
- As emergency back you can use a Docker container from Duke OIT
 1. Go to <https://cmgr.oit.duke.edu/> and login
 2. Select [Reserve a Container](#) and find a container for any Sta course
 3. Click the link under my reservations to create your environment

Local R + RStudio

If working locally you should make sure that your environment meets the following requirements:

- latest R (4.5.2)
- latest RStudio (2026.01.0)
- working git installation
- ability to create ssh keys (for GitHub authentication)
- *All* R packages updated to their latest version from CRAN

Support policy for local installs - we will try to help you troubleshoot if we can but reserve the right to tell you to use the dept server.

GitHub

- We will be using a GitHub organization for this course github.com/sta323-sp26
- All assignments will be distributed and collected via GitHub
- All of your work and your membership (enrollment) in the organization is private
- We will be distributing a survey this weekend to collection your GitHub account names
 - Before lab you will be invited to the course organization.
- All course related repositories will be created for you

Before Monday

- Create a GitHub account if you don't have one
- Complete the course survey
- Make sure you can login in to the Department's RStudio server
<https://rstudio.stat.duke.edu>
- Setup ssh key authentication with GitHub, see
https://github.com/DukeStatSci/github_auth_guide

**In \mathbb{R} (almost)
everything is a vector**

Vectors

The fundamental building block of data in R are vectors (collections of related values, objects, etc).

R has two types of vectors:

- atomic vectors (*vectors*)
 - homogeneous collections of the *same* type (e.g. all `true/false` values, all numbers, or all character strings).
- generic vectors (*lists*)
 - heterogeneous collections of *any* type of R object, even other lists (meaning they can have a hierarchical/tree-like structure).

Atomic Vectors

Atomic Vectors

R has six atomic vector types.

<code>typeof()</code>	<code>mode()</code>
logical	logical
double	numeric
integer	numeric
character	character
complex	complex
raw	raw

We can check the type of any object in R using the `typeof()` function. `mode()` is a higher level abstraction used to group similar types together.

There are additional types in R, e.g. `list`, `closure`, `environment`, etc. We will see these in the next couple of lectures. See

Logical - boolean values (TRUE and FALSE)

```
1 typeof(TRUE)
```

```
[1] "logical"
```

```
1 typeof(FALSE)
```

```
[1] "logical"
```

```
1 mode(TRUE)
```

```
[1] "logical"
```

```
1 mode(FALSE)
```

```
[1] "logical"
```

R will let you use **T** and **F** as shortcuts to **TRUE** and **FALSE**, this is bad practice as these values are actually **global variables** that can be overwritten.

```
1 T
```

```
[1] TRUE
```

```
1 T = "FALSE"
```

```
2 T
```

```
[1] "FALSE"
```

character - text strings

Either single or double quotes are fine, the opening and closing quote must match.

```
1 typeof("hello")
```

```
[1] "character"
```

```
1 typeof('world')
```

```
[1] "character"
```

```
1 mode("hello")
```

```
[1] "character"
```

```
1 mode('world')
```

```
[1] "character"
```

Quote characters can be included by escaping or using a non-matching quote.

```
1 "abc '123"
```

```
[1] "abc '123"
```

```
1 'abc"123'
```

```
[1] "abc\"123"
```

```
1 "abc\"123"
```

```
[1] "abc\"123"
```

```
1 'abc\'123'
```

```
[1] "abc '123"
```

Numeric types

`double` - floating point values (these are the default numerical type)

```
1 typeof(1.33)
```

```
[1] "double"
```

```
1 typeof(7)
```

```
[1] "double"
```

```
1 mode(1.33)
```

```
[1] "numeric"
```

```
1 mode(7)
```

```
[1] "numeric"
```

`integer` - integer values (literals are indicated by the `L` suffix)

```
1 typeof( 7L )
```

```
[1] "integer"
```

```
1 typeof( 1:3 )
```

```
[1] "integer"
```

```
1 mode( 7L )
```

```
[1] "numeric"
```

```
1 mode( 1:3 )
```

```
[1] "numeric"
```

Combining / Concatenation

Atomic vectors can be constructed using the combine `c()` function.

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

```
[1] 1 2 3
```

```
1 c("Hello", "World!")
```

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

```
1 c(1, 1:10)
```

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

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

```
[1] 1 2 3
```

Type predicates

- `is.logical(x)` - returns `TRUE` if `x` has *type* `logical`.
- `is.character(x)` - returns `TRUE` if `x` has *type* `character`.
- `is.double(x)` - returns `TRUE` if `x` has *type* `double`.
- `is.integer(x)` - returns `TRUE` if `x` has *type* `integer`.
- `is.numeric(x)` - returns `TRUE` if `x` has *mode* `numeric`.

```
1 is.integer(1)
```

```
[1] FALSE
```

```
1 is.integer(1L)
```

```
[1] TRUE
```

```
1 is.integer(3:7)
```

```
[1] TRUE
```

```
1 is.double(1)
```

```
[1] TRUE
```

```
1 is.double(1L)
```

```
[1] FALSE
```

```
1 is.double(3:8)
```

```
[1] FALSE
```

```
1 is.numeric(1)
```

```
[1] TRUE
```

```
1 is.numeric(1L)
```

```
[1] TRUE
```

```
1 is.numeric(3:7)
```

```
[1] TRUE
```

Other useful predicates

- `is.atomic(x)` - returns `TRUE` if `x` is an *atomic vector*.
- `is.list(x)` - returns `TRUE` if `x` is a *list* (generic vector).
- `is.vector(x)` - returns `TRUE` if `x` is either an *atomic* or *generic* vector.

```
1 is.atomic(c(1,2,3))
```

```
[1] TRUE
```

```
1 is.list(c(1,2,3))
```

```
[1] FALSE
```

```
1 is.vector(c(1,2,3))
```

```
[1] TRUE
```

```
1 is.atomic(list(1,2,3))
```

```
[1] FALSE
```

```
1 is.list(list(1,2,3))
```

```
[1] TRUE
```

```
1 is.vector(list(1,2,3))
```

```
[1] TRUE
```

Type Coercion

R is a dynamically typed language – it will automatically convert between most types without raising warnings or errors. Keep in mind that atomic vectors must always contain values of the same type.

```
1 c(1, "Hello")
```

```
[1] "1"      "Hello"
```

```
1 c(FALSE, 3L)
```

```
[1] 0 3
```

```
1 c(1.2, 3L)
```

```
[1] 1.2 3.0
```

```
1 c(FALSE, "Hello")
```

```
[1] "FALSE" "Hello"
```

Operator coercion

Builtin operators and functions (e.g. `+`, `&`, `log()`, etc.) will generally attempt to coerce values to an appropriate type for the given operation (numeric for math, logical for logical, etc.)

```
1 3.1+1L
```

```
[1] 4.1
```

```
1 5 + FALSE
```

```
[1] 5
```

```
1 TRUE & FALSE
```

```
[1] FALSE
```

```
1 TRUE & 7
```

```
[1] TRUE
```

```
1 log(1)
```

```
[1] 0
```

```
1 log(TRUE)
```

```
[1] 0
```

```
1 TRUE | FALSE
```

```
[1] TRUE
```

```
1 FALSE | !5
```

```
[1] FALSE
```

Explicit coercion

Most of the `is` functions we just saw have an `as` variant which can be used for *explicit* coercion.

```
1 as.logical(5.2)
```

```
[1] TRUE
```

```
1 as.character(TRUE)
```

```
[1] "TRUE"
```

```
1 as.integer(pi)
```

```
[1] 3
```

```
1 as.numeric(FALSE)
```

```
[1] 0
```

```
1 as.double("7.2")
```

```
[1] 7.2
```

```
1 as.double("one")
```

```
Warning: NAs introduced by coercion
```

```
[1] NA
```

Missing Values

Missing Values

R uses `NA` to represent missing values in its data structures.

What may not be obvious is that there are different `NA`s for the different atomic types.

<pre>1 NA</pre>	<pre>1 typeof(NA)</pre>	<pre>1 typeof(NA_character_)</pre>
<pre>[1] NA</pre>	<pre>[1] "logical"</pre>	<pre>[1] "character"</pre>
<pre>1 NA+1</pre>	<pre>1 typeof(NA+1)</pre>	<pre>1 typeof(NA_real_)</pre>
<pre>[1] NA</pre>	<pre>[1] "double"</pre>	<pre>[1] "double"</pre>
<pre>1 NA+1L</pre>	<pre>1 typeof(NA+1L)</pre>	<pre>1 typeof(NA_integer_)</pre>
<pre>[1] NA</pre>	<pre>[1] "integer"</pre>	<pre>[1] "integer"</pre>
<pre>1 c(NA, "")</pre>	<pre>1 typeof(c(NA, ""))</pre>	<pre>1 typeof(NA_complex_)</pre>
<pre>[1] NA ""</pre>	<pre>[1] "character"</pre>	<pre>[1] "complex"</pre>

NA stickiness

As **NA**s represent missing values (most) calculations using them return a missing value.

```
1 1 + NA
```

```
[1] NA
```

```
1 1 / NA
```

```
[1] NA
```

```
1 NA * 5
```

```
[1] NA
```

```
1 sqrt(NA)
```

```
[1] NA
```

```
1 3^NA
```

```
[1] NA
```

```
1 sum(c(1, 2, 3, NA))
```

```
[1] NA
```

Aggregation / summarization functions (e.g. `sum()`, `mean()`, `sd()`, etc.) will often have an `na.rm` argument which *removes* the missing values from the calculation.

```
1 sum(c(1, 2, 3, NA), na.rm = TRUE)
```

```
[1] 6
```

```
1 mean(c(1, 2, 3, NA), na.rm = TRUE)
```

```
[1] 2
```

NAs are not always sticky

A useful mental model for **NAs** is to consider them as a unknown value that could take any of the possible values for a given type.

For numbers or characters this isn't helpful, but for a logical value it must either be **TRUE** or **FALSE** which is relevant for certain calculations.

```
1 TRUE & NA
```

```
[1] NA
```

```
1 FALSE & NA
```

```
[1] FALSE
```

```
1 TRUE | NA
```

```
[1] TRUE
```

```
1 FALSE | NA
```

```
[1] NA
```

Other Special values (double)

These are defined as part of the IEEE floating point standard (not unique to R)

- **NaN** - Not a number
- **Inf** - Positive infinity
- **-Inf** - Negative infinity

1 **-1** / **0**

[1] **-Inf**

1 **0** / **0**

[1] **NaN**

1 **1/0** + **1/0**

[1] **Inf**

1 **Inf** - **Inf**

[1] **NaN**

1 **NaN** / **NA**

[1] **NA**

1 **NaN** * **NA**

[1] **NA**

Testing for Inf and NaN

there are predicate functions for testing for these types of values

```
1 is.finite(Inf)
```

```
[1] FALSE
```

```
1 is.infinite(-Inf)
```

```
[1] TRUE
```

```
1 is.nan(Inf)
```

```
[1] FALSE
```

```
1 Inf > 1
```

```
[1] TRUE
```

```
1 is.finite(NaN)
```

```
[1] FALSE
```

```
1 is.infinite(NaN)
```

```
[1] FALSE
```

```
1 is.nan(NaN)
```

```
[1] TRUE
```

```
1 -Inf > 1
```

```
[1] FALSE
```

```
1 is.finite(NA)
```

```
[1] FALSE
```

```
1 is.infinite(NA)
```

```
[1] FALSE
```

```
1 is.nan(NA)
```

```
[1] FALSE
```

Coercion for infinity and NaN

First remember that `Inf`, `-Inf`, and `NaN` are doubles, however their coercion behavior is not the same as other doubles

```
1 as.integer(Inf)
```

Warning: NAs introduced by coercion to integer range

```
[1] NA
```

```
1 as.integer(NaN)
```

```
[1] NA
```

```
1 as.logical(Inf)
```

```
[1] TRUE
```

```
1 as.logical(-Inf)
```

```
[1] TRUE
```

```
1 as.logical(NaN)
```

```
[1] NA
```

```
1 as.character(Inf)
```

```
[1] "Inf"
```

```
1 as.character(-Inf)
```

```
[1] "-Inf"
```

```
1 as.character(NaN)
```

```
[1] "NaN"
```

Exercise 1

Part 1

What is the type of the following vectors? Explain why they have that type.

- 1 `c(1, NA+1L, "C")`
- 2 `c(1L / 0, NA)`
- 3 `c(1:3, 5)`
- 4 `c(3L, NaN+1L)`
- 5 `c(NA, TRUE)`

Part 2

Considering only the four (common) data types, what is R's implicit type conversion hierarchy (from highest priority to lowest priority)?