



# **Coding Basics and Importing Data**

Joseph Rudolf  
AACC 2021

September 26, 2021	Session	Instructor
8:30 am – 8:45 am	Course Introduction	Patrick Mathias
8:45 am – 9:30 am	Intro to R and Reproducible Reporting	Joseph Rudolf
9:45 am – 10:30 am	Coding Basics and Importing Data	Joseph Rudolf
10:45 am – 11:30 am	Data Visualization	Patrick Mathias
LUNCH		
12:30 pm – 1:30 pm	Data Transformation	Patrick Mathias
1:45 pm – 2:45 pm	Grouping and Summarizing Data	Joseph Rudolf
3:00 pm – 3:30 pm	Dashboard Demo and Course Wrap Up	Patrick Mathias

# Lesson Goals

1. Learn some fundamental of coding

# Lesson Objectives

1. Define and use functions
2. Define and create objects in the environment
3. Install and load packages
4. Import data into R
5. Interact with a dataframe



# The Basics of Coding

# The Basics of Coding: Calculation

- R is a calculator!



A screenshot of an R console window. The window has a light gray background. On the left, there is a vertical line of numbers 1 through 6, with a small triangle pointing to the right next to each number. The text in the console is as follows: line 1 is empty; line 2 is ````{r}```; line 3 is empty; line 4 is empty; line 5 is empty; line 6 is ````. To the right of the text, there are three icons: a gear icon, a green minus sign icon, and a green right-pointing triangle icon. Below the main text area, there is a small window showing the output: ``[1] 7``. To the right of the output, there are three icons: a document icon, a green up-pointing triangle icon, and a red X icon.

```
1  
2 ````{r}```  
3  
4  
5  
6 ````  
[1] 7
```

press play  
button to execute  
code

answer returned  
here

# The Basics of Coding: Functions

- Code that extends our reach beyond the basic operators

```
1  
2  ```{r}  
3  
4  abs(-77)  
5  
6  ```  
[1] 77
```

function  
(does stuff)

argument  
(input)

abs(-77)

# Putting Functions to Work

- We can use functions to do more than simple math, we can make things!
- We can create a series of integers (a vector) using the `seq()` function

```
1  
2 ``{r}  
3  
4 seq(from=5, to=150, by=10)  
5  
6 ```
```

```
[1]  5 15 25 35 45 55 65 75 85 95 105 115 125 135 145
```

# The Basics of Coding: Objects

- Objects are the container for your output

object  
(stores output)

function  
(does stuff)

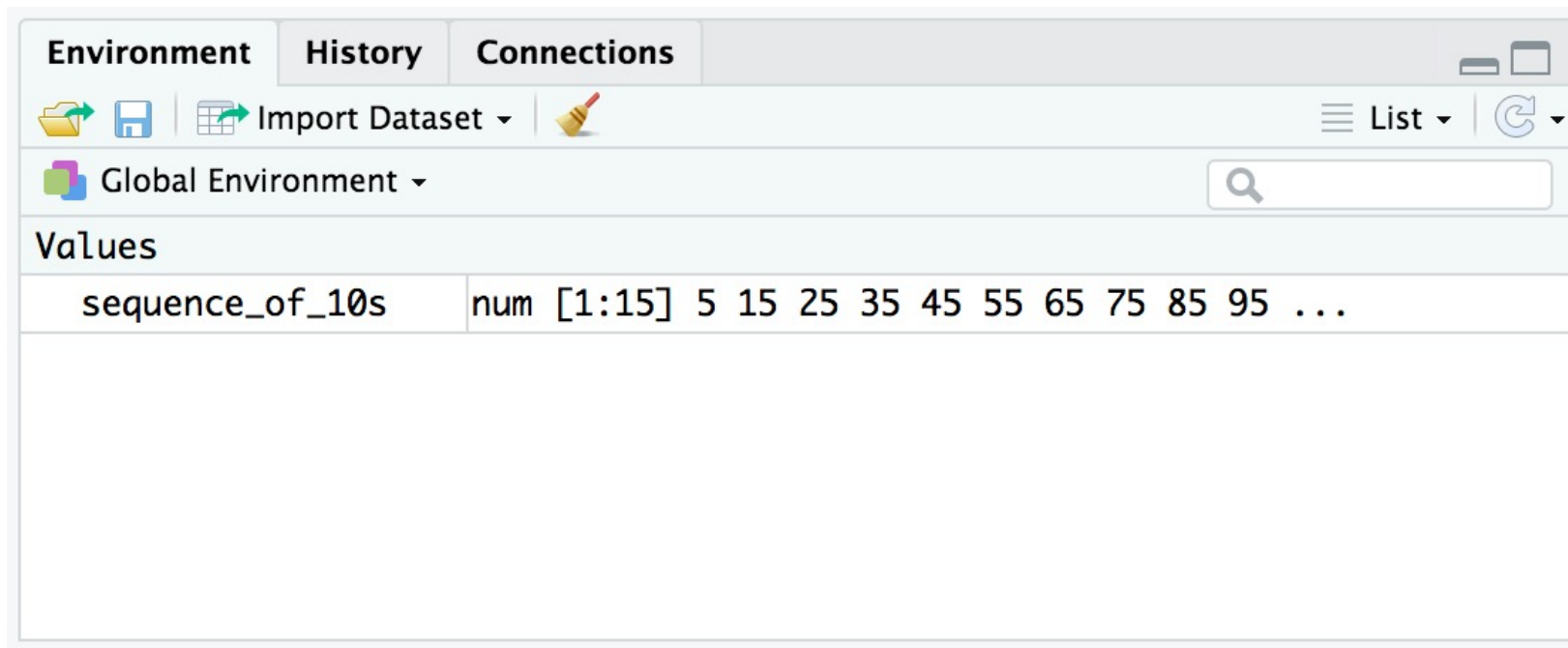
arguments  
(input)

```
sequence_of_10s <- seq(from=5, to=150, by=10 )
```



# Checking the contents of an object

- The environment tab shows us the objects we have created.



# Bending objects to your will

- Once we have created an object we can start to interact with it.
- This includes passing our objects to other functions... Whoa!

```
1  
2  ```{r}  
3  
4  min(sequence_of_10s)  
5  
6  ```
```

[1] 5

```
1  
2  ```{r}  
3  
4  max(sequence_of_10s)  
5  
6  ```
```

[1] 145

# Your Turn #1

I've written some code to create a sequence from 0 to 500 in increments of 25 called `sequence_of_25s`. Ultimately I want to calculate the median value of this sequence. Unfortunately I've made some mistakes in my code and I am hoping you can help me find them.

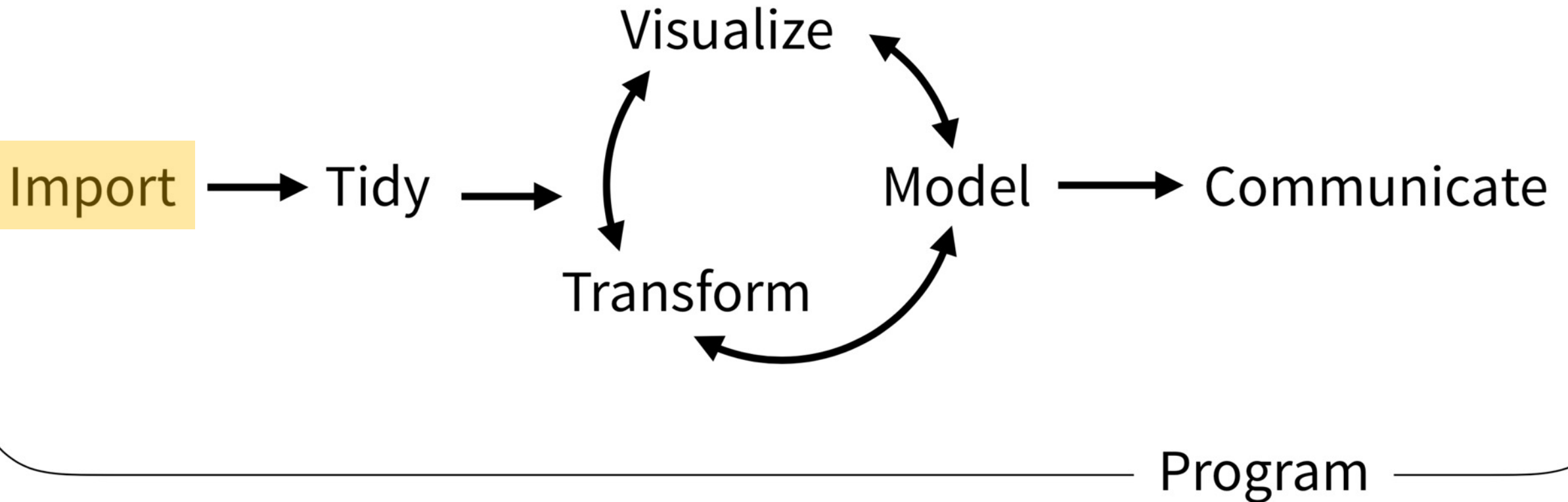
```
1  
2  ```{r}  
3  
4  sequence_of_25s -< seq(from=0 to=50, by=25)  
5  
6  ```  
7  
8  ```{r}  
9  
10 median(sequence of_25s]  
11  
12 ```  
13
```

01:00



# Importing Data

# The Data Analysis Pipeline



plain text  
("flat") file



header row

02-example		
Name	MRN	DOB
Santa Claus	12345	1/1/01
Roger Rabbit	67890	12/12/69
Kermit the Frog	24680	2/2/22

rectangular  
structure

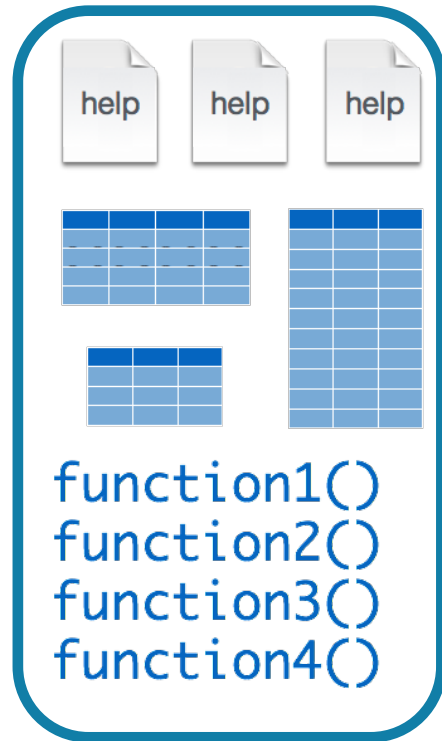
# Tidyverse: R Packages for Data Science

- A consistent way to organize data
- Human readable, concise, consistent code
- Build pipelines from atomic data analysis steps



# Installing and loading R packages

tidyverse



```
install.packages("tidyverse")
```

Downloads files to computer

**1 x per computer**

```
library("tidyverse")
```

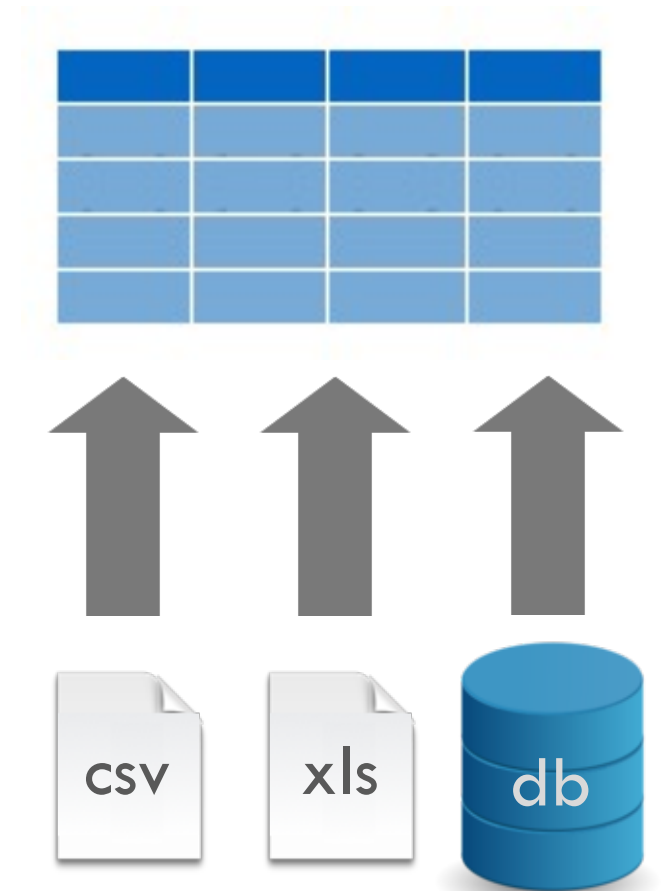
Loads package

**1 x per R Session**



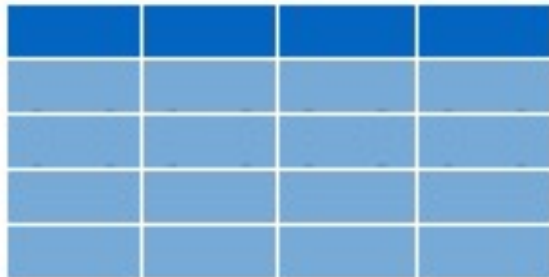
# Dataframes: Beyond the Vector

- Dataframe is the term for a table
- Dataframes are composed:  
Columns (Variables)  
Rows (Observations)
- Dataframes are objects and can be acted on like other objects



# read\_csv()

```
data_frame <- read_csv(file_name)
```







function  
(does stuff)

```
data_frame <- read_csv(file_name)
```

function  
(does stuff)

argument  
(input)

```
data_frame <- read_csv(file_name)
```

object  
(stores output)

function  
(does stuff)

argument  
(input)

```
data_frame <- read_csv(file_name)
```

object  
(stores output)

function  
(does stuff)

argument  
(input)

```
data_frame <- read_csv(file_name)
```

assignment operator  
("gets")

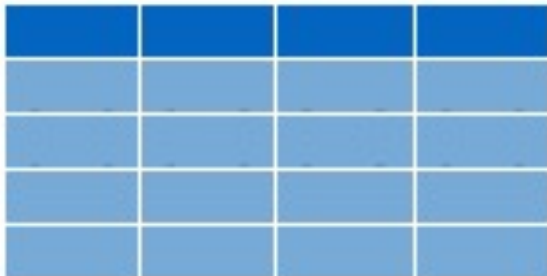
# read\_csv()

data frame to  
read data into

name of  
CSV file

```
covid_testing <- read_csv("covid_testing.csv")
```

covid\_testing







covid\_testing.csv



# Your Turn #2

In the MISC pane, select the folder:  
“exercises”

Select the R Markdown file:  
“02 – Importing and Exploring Data.Rmd”

In the Editor pane, follow the instructions to complete the exercise.

05:00



# Recap

**Functions** do stuff. They accept **Arguments** to define parameters. We can store the output of functions in **Objects** using the assignment operator ( `<-` ).

**Packages** extend the functionality of R. They need to be installed once per computer and loaded each session.

**Importing Data** is the first step data analysis pipeline. `read_csv()` is a function from the tidyverse that we can use for importing data.



**What else?**



# Data Import :: CHEAT SHEET

R's **tidyverse** is built around **tidy data** stored in **tibbles**, which are enhanced data frames.



The front side of this sheet shows how to read text files into R with **readr**.



The reverse side shows how to create tibbles with **tibble** and to layout tidy data with **tidyr**.

## OTHER TYPES OF DATA

Try one of the following packages to import other types of files

- **haven** - SPSS, Stata, and SAS files
- **readxl** - excel files (.xls and .xlsx)
- **DBI** - databases
- **jsonlite** - json
- **xml2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)

## Save Data

Save **x**, an R object, to **path**, a file path, as:

### Comma delimited file

```
write_csv(x, path, na = "NA", append = FALSE,
  col_names = !append)
```

### File with arbitrary delimiter

```
write_delim(x, path, delim = " ", na = "NA",
  append = FALSE, col_names = !append)
```

### CSV for excel

```
write_excel_csv(x, path, na = "NA", append =
  FALSE, col_names = !append)
```

### String to file

```
write_file(x, path, append = FALSE)
```

### String vector to file, one element per line

```
write_lines(x, path, na = "NA", append = FALSE)
```

### Object to RDS file

```
write_rds(x, path, compress = c("none", "gz",
  "bz2", "xz"), ...)
```

### Tab delimited files

```
write_tsv(x, path, na = "NA", append = FALSE,
  col_names = !append)
```

## Read Tabular Data

- These functions share the common arguments:

```
read_(file, col_names = TRUE, col_types = NULL, locale = default_locale(), na = c("", "NA"),
  quoted_na = TRUE, comment = "", trim_ws = TRUE, skip = 0, n_max = Inf, guess_max = min(1000,
  n_max), progress = interactive())
```

a,b,c  
1,2,3  
4,5,NA

A	B	C
1	2	3
4	5	NA

### Comma Delimited Files

**read\_csv("file.csv")**

To make file.csv run:

```
write_file(x = "a,b,c\n1,2,3\n4,5,NA", path = "file.csv")
```

a;b;c  
1;2;3  
4;5;NA

A	B	C
1	2	3
4	5	NA

### Semi-colon Delimited Files

**read\_csv2("file2.csv")**

```
write_file(x = "a;b;c\n1;2;3\n4;5;NA", path = "file2.csv")
```

a|b|c  
1|2|3  
4|5|NA

A	B	C
1	2	3
4	5	NA

### Files with Any Delimiter

**read\_delim("file.txt", delim = "|")**

```
write_file(x = "a|b|c\n1|2|3\n4|5|NA", path = "file.txt")
```

a b c  
1 2 3  
4 5 NA

A	B	C
1	2	3
4	5	NA

### Fixed Width Files

**read\_fwf("file.fwf", col\_positions = c(1, 3, 5))**

```
write_file(x = "a b c\n1 2 3\n4 5 NA", path = "file.fwf")
```

### Tab Delimited Files

**read\_tsv("file.tsv")** Also **read\_table()**.

```
write_file(x = "a\tb\tc\n1\t2\t3\n4\t5\tNA", path = "file.tsv")
```

## USEFUL ARGUMENTS

a,b,c  
1,2,3  
4,5,NA

### Example file

```
write_file("a,b,c\n1,2,3\n4,5,NA","file.csv")  
f<- "file.csv"
```

1	2	3
4	5	NA

### Skip lines

```
read_csv(f, skip = 1)
```

A	B	C
1	2	3
4	5	NA

### No header

```
read_csv(f, col_names = FALSE)
```

A	B	C
1	2	3

### Read in a subset

```
read_csv(f, n_max = 1)
```

x	y	z
1	2	3
4	5	NA

### Provide header

```
read_csv(f, col_names = c("x", "y", "z"))
```

A	B	C
NA	2	3
4	5	NA

### Missing Values

```
read_csv(f, na = c("1", ""))
```

## Read Non-Tabular Data

### Read a file into a single string

```
read_file(file, locale = default_locale())
```

### Read each line into its own string

```
read_lines(file, skip = 0, n_max = -1L, na = character(),
  locale = default_locale(), progress = interactive())
```

### Read Apache style log files

```
read_log(file, col_names = FALSE, col_types = NULL, skip = 0, n_max = -1, progress = interactive())
```

### Read a file into a raw vector

```
read_file_raw(file)
```

### Read each line into a raw vector

```
read_lines_raw(file, skip = 0, n_max = -1L,
  progress = interactive())
```

## Data types

readr functions guess the types of each column and convert types when appropriate (but will NOT convert strings to factors automatically).

A message shows the type of each column in the result.

```
## Parsed with column specification:
## cols(
##   age = col_integer(),
##   sex = col_character(),
##   earn = col_double()
## )
```

age is an integer  
sex is a character  
earn is a double (numeric)

1. Use **problems()** to diagnose problems

```
x <- read_csv("file.csv"); problems(x)
```

2. Use a **col\_** function to guide parsing

- **col\_guess()** - the default
- **col\_character()**
- **col\_double()**, **col\_euro\_double()**
- **col\_datetime()** (format = "") Also **col\_date** (format = ""), **col\_time** (format = "")
- **col\_factor** (levels, ordered = FALSE)
- **col\_integer()**
- **col\_logical()**
- **col\_number()**, **col\_numeric()**
- **col\_skip()**

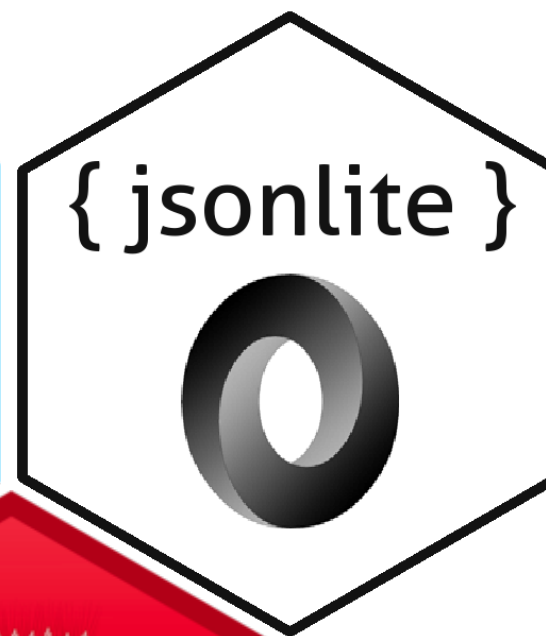
```
x <- read_csv("file.csv", col_types = cols(
  A = col_double(),
  B = col_logical(),
  C = col_factor()))
```

3. Else, read in as character vectors then parse with a **parse\_** function.

- **parse\_guess()**
- **parse\_character()**
- **parse\_datetime()** Also **parse\_date()** and **parse\_time()**
- **parse\_double()**
- **parse\_factor()**
- **parse\_integer()**
- **parse\_logical()**
- **parse\_number()**

```
x$A <- parse_number(x$A)
```





# Databases



[Microsoft SQL Serve](#)



[MonetDB](#)



[MongoDB](#)



[MySQL](#)



[Netezza](#)



[Oracle](#)



[Amazon Redshift](#)



[Apache Hive](#)



[Apache Impala](#)



[Athena](#)



[Cassandra](#)



[Google BigQuery](#)



[Other Databases](#)



[PostgreSQL](#)



[SQLite](#)



[Salesforce](#)



[Teradata](#)

# R Interface to Python



```
```{python}  
import pandas  
covid_testing.info()  
```
```