Data Import

with readr, tibble, and tidyr

Cheat Sheet



R's tidyverse is built around tidy data stored in **:ibbles**, an enhanced version of a data frame.



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)

Write functions

Save x, an R object, to path, a file path, with:

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

Tibble/df to comma delimited file.

write_delim(x, path, delim = " ", na = "NA", append = FALSE, col_names = !append) Tibble/df to file with any delimiter.

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

Tibble/df to a CSV for excel

write_file(x, path, append = FALSE) String to file.

String vector to file, one element per line. write_lines(x, path, na = "NA", append =

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

Object to RDS file.

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

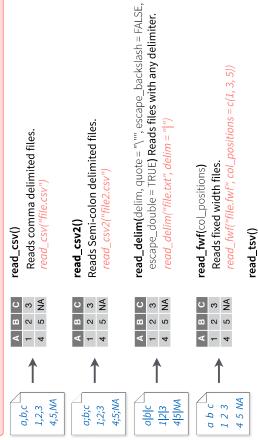
Fibble/df to tab delimited files.

Read functions

These functions share the common arguments:

Read tabular data to tibbles

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())



Useful arguments

x = read_csv("a,b,c\n1,2,3\n4,5,NA")) write_csv (path = "file.csv", Example file 4,5,NA a,b,c1,2,3

A B C

က

1

col_names = FALSE)

read_csv("file.csv'

No header

1 A B C

4 5 NA

read_csv("file.csv",

skip = 1

Skip lines

read_csv("file.csv", Read in a subset $n_max = 1$

A B NA NA NA က 7

 $col_names = c("x", "y", "z"))$

ღ ¥

read_csv("file.csv",

Provide header

x × ပ

ш N

na = c("4", "5", "")) read_csv("file.csv", **Missing Values**

Read non-tabular data

read_file, locale = default_locale()) Read a file into a single string.

read_file_raw(file)

read_lines(file, skip = 0, n_max = -1L, locale = default_locale(), na = character(), progress = Read a file into a raw vector.

Read each line into its own string.

interactive())

read_lines_raw(file, skip = 0, n_max = -1L, Read each line into a raw vector. progress = interactive())

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

Apache style log files.

Parsing 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()
##) characte 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 = "") and col_time(format = "")
 - col_factor(levels, ordered = FALSE)
 - col_integer()
 - col_logical()
- col_number()

Reads tab delimited files. Also read_table()

read_tsv("file.tsv")

- col_numeric() col_skip()
- x <- read_csv("file.csv", col_types = cols(() = col_double()

B = col_logical(), $C = col_factor()$ 3. Else, read in as character vectors then parse with a parse_function.

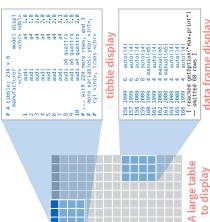
- parse_guess(x, na = c("", "NA"), locale = default_locale())
- parse_character(x, na = c("", "NA"), locale = default_locale())
- **parse_datetime(**x, format = "", na = c("", "NA"),
 - **parse_double(**x, na = c("", "NA"), locale = and parse_time()
- parse_factor(x, levels, ordered = FALSE, na = c("", "NA"), locale = default_locale()) default_locale())
- **parse_integer(**x, na = c("", "NA"), locale = parse_logical(x, na = c("", "NA"), locale = default_locale())
- parse_number(x, na = c("", "NA"), locale = default_locale())
 - default_locale())

κ\$A <- parse_number(x\$A)

Tibbles - an enhanced data frame

storing tabular data, the tibble. Tibbles inherit the The tibble package provides a new S3 class for data frame class, but improve two behaviors:

- Display When you print a tibble, R provides a concise view of the data that fits on one screen.
- . Subsetting [always returns a new tibble, [[and \$ always return a vector.
- No partial matching You must use full column names when subsetting



- Control the default appearance with options: **options**(tibble.print_max = n,
 - tibble.print_min = m, tibble.width = Inf) View entire data set with View(x, title) or **glimpse(**x, width = NULL, ...
 - Revert to data frame with **as.data.frame()** (required for some older packages)

Construct a tibble in two ways

A tibble: 3 × 2 × y <int> <dbl> **Both make** this tibble Construct by columns. Construct by rows. tibble(x =tribble(...) tibble(...) tribble(

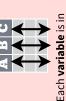
as_tibble(x, ...) Convert data frame to tibble.

enframe(x, name = "name", value = "value") Converts named vector to a tibble with a

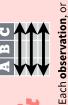
names column and a values column. is_tibble(x) Test whether x is a tibble.

Tidy data is a way to organize tabular data. It provides a consistent data structure across packages. Fidy data: A table is tidy if:

Tidv Data with tidvr











Makes variables easy



^ B * * **4**

Preserves cases during

vectorized operations

to access as vectors

case, is in its own row

ts own column

Reshape Data - change the layout of values in a table

Jse gather() and spread() to reorganize the values of a table into a new layout. Each uses the idea of a key column: value column pair.

gather(data, key, value, ..., na.rm = FALSE, convert = FALSE, factor_key = FALSE)

column, gathering the column values into a Gather moves column names into a key single value column.

into the column names, spreading the values of a value column across the new columns that result.

Spread moves the unique values of a key column

spread(data, key, value, fill = NA, convert = FALSE,

drop = TRUE, sep = NULL

country 1999 2000 A 0.7K 2K 212K 213K 37K table4a A W O



country year type count key value

country year A 1999 A 1999 cases 212K 2000 cases 213K 2000 cases 80K 1999 cases 2000 pop 2000 cases 1999 pop 2000 pop 1999 pop 1999 pop 2000 A B B В В

2000 80K 1999 212K

1999 37K

2000

spread(table2, type, count) key value

gather(table4a, `1999`, `2000` key = "year", value = "cases")

Handle Missing Values

Drop rows containing NA's in ... columns. drop_na(data, ...)

fill(data, ..., .direction = c("down", "up")**)** Fill in NA's in ... columns with most recent non-NA values.

x1 x2

x1 x2

E D C B A x1 x2

x1 x2 x1 x2

replace na(x, list(x2 = 2), x2)

fill(x, x2)

drop_na(x, x2)

unite(data, col, ..., sep = "_", remove = TRUE)

Collapse cells across several columns to make a single column.

country year Afghan Afghan Brazil Brazil China country century y Afghan 19 Afghan 20 Brazil China Brazi

unite(table5, century, year,

Expand Tables - quickly create tables with combinations of values

complete(data, ..., fill = list())

Adds to the data missing combinations of the values of the variables listed in ... complete(mtcars, cyl, gear, carb)

expand(data, ...

Create new tibble with all possible combinations of the values of the variables listed in ... expand(mtcars, cyl, gear, carb)

Split and Combine Cells

Use these functions to split or combine cells into ndividual, isolated values.

separate(data, col, into, sep = "[^[:alnum:]]+", remove = TRUE, convert = FALSE, extra = "warn", fill = "warn", ... Separate each cell in a column to make several columns.

dod	19M	20M	172	174	Ļ	Ħ		
cases	0.7K	2K	37K	80K	212K	213K		
year	1999	2000	1999	2000	1999	2000		
country	∢	∢	В	В	O	O		
1								
rate	0.7K/19M	2K/20M	37K/172M	80K/174M	212K/1T	213K/1T		
year	1999	2000	1999	2000	1999	2000		
country	∢	∢	В	മ	ပ	ပ		

separate_rows(table3, rate, into = c("cases", "pop")) **separate_rows**(data, ..., sep = "[^[:alnum:].]+", convert = FALSE Separate each cell in a column to make several rows. Also separate_rows_()

	0.7K	19M	2K	20M	37K	172M	80K	174M	212K		
year	1999	1999	2000	2000	1999	1999	2000	2000	1999		
country	∢	∢	∢	∢	В	В	В	В	O		
↑											
rate	0.7K/19M	2K/20M	37K/172M	80K/174M	212K/1T	213K/1T					
year	1999	2000	1999	2000	1999	2000					
country	∢	∢	В	В	ပ	ပ					

separate_rows(table3, rate)

Replace NA's by column.

replace_na(_{data,}

replace = list(),

2000 213K

00

1999

O O col = "year", sep = ""