# Reading and Writing Data

## readr and haven

2024-07-18

# readr

| Function | Reads |
|---|---|
| `read_csv()` | Comma separated values |
| `read_csv2()` | Semi-colon separate values |
| `read_delim()` | General delimited files |
| `read_fwf()` | Fixed width files |
| `read_log()` | Apache log files |
| `read_table()` | Space separated files |
| `read_tsv()` | Tab delimited values |

# Importing Data

```
1  dataset <- read_csv("file_name.csv")
2  dataset
```

# R functions

```r
x <- f(arg = 1)
```

# R functions

x <- f(arg = 1)

function
name

arguments

# R functions

*this saves it in your global environment*

x <- f(arg = 1)

assign
results of
f() to x

the name of
your results

# *Your Turn 1*

## Find **diabetes.csv** on your computer. Then read it into an object. Then view the results.

...

```
1  diabetes <- read_csv("diabetes.csv")
```

```
1  diabetes
```

# A tibble: 403 × 19
```
        id    chol stab.glu    hdl  ratio  glyhb location       age
     <dbl>  <dbl>    <dbl>   <dbl>  <dbl>  <dbl> <chr>         <dbl>
 1   1000    203       82      56   3.60   4.31 Buckingham       46
 2   1001    165       97      24   6.90   4.44 Buckingham       29
 3   1002    228       92      37   6.20   4.64 Buckingham       58
 4   1003     78       93      12   6.5    4.63 Buckingham       67
 5   1005    249       90      28   8.90   7.72 Buckingham       64
 6   1008    248       94      69   3.60   4.81 Buckingham       34
 7   1011    195       92      41   4.80   4.84 Buckingham       30
 8   1015    227       75      44   5.20   3.94 Buckingham       37
 9   1016    177       87      49   3.60   4.84 Buckingham       45
10   1022    263       89      40   6.60   5.78 Buckingham       55
```

# Tibbles

`data.frames` are the basic form of rectangular data in R (columns of variables, rows of observations)

`read_csv()` reads the data into a `tibble`, a modern version of the data frame.

a tibble **is** a data frame

# Missing values

It's common to use codes for **missing values** (-99, 8888)

The na option can change these values to NA

```
1  read_csv(
2    "a,b,c,d
3    1,-99,3,4
4    5,6,-99,8",
5    na = "-99"
6  )
```

```
# A tibble: 2 × 4
      a     b     c     d
  <dbl> <dbl> <dbl> <dbl>
1     1    NA     3     4
2     5     6    NA     8
```

# Parsing data types

The read functions in readr try to *guess* each data type, but sometimes it's *wrong*

To tell readr how to parse the columns, add the argument **col_types** to read_csv()

```r
1  diabetes <- read_csv(
2    "diabetes.csv",
3    col_types = list(id = col_character())
4  )
```

# Parsing data types

**Or use a string for each variable type:** `col_type = "cci"`

# Parsing data types

## Or use a string for each variable type: col_type = "cci"

| letter | type |
| --- | --- |
| c | character |
| i | integer |
| n | number |
| d | double |
| l | logical |
| D | date |
| T | date time |
| t | time |
| ? | guess the type |
| _ or - | skip the column |

# *Your Turn 2*

## Set the 4 column types to be: integer, double, character, and unknown (guess)

```
1  read_csv(
2    "a,b,c,d
3    1,2,3,4
4    5,6,7,8",
5    col_types = ""
6  )
```

# *Your Turn 2*

## Set the 4 column types to be: integer, double, character, and unknown (guess)

```
1  read_csv(
2    "a,b,c,d
3    1,2,3,4
4    5,6,7,8",
5    col_types = "idc?"
6  )
```

```
# A tibble: 2 × 4
      a     b c         d
  <int> <dbl> <chr> <dbl>
1     1     2 3         4
2     5     6 7         8
```

# haven

| Function | Software |
|---|---|
| read_sas() | SAS |
| read_xpt() | SAS |
| read_spss() | SPSS |
| read_sav() | SPSS |
| read_por() | SPSS |
| read_stata() | Stata |
| read_dta() | Stata |

# Heads up!

haven is *not* a core member of the tidyverse. That means you need to load it with `library(haven)`.

# *Your Turn 3*

**There are several versions of the diabetes file besides CSV. Pick a file format you or your colleagues use and import them using the corresponding function from haven.**

# *Your Turn 3*

```r
library(haven)
diabetes <- read_sas("diabetes.sas7bdat")
```

# Your Turn 3

```
1  diabetes
```

```
# A tibble: 403 × 19
      id  chol stab_glu   hdl ratio glyhb location       age
   <dbl> <dbl>    <dbl> <dbl> <dbl> <dbl> <chr>        <dbl>
 1  1000   203       82    56  3.60  4.31 Buckingham      46
 2  1001   165       97    24  6.90  4.44 Buckingham      29
 3  1002   228       92    37  6.20  4.64 Buckingham      58
 4  1003    78       93    12  6.5   4.63 Buckingham      67
 5  1005   249       90    28  8.90  7.72 Buckingham      64
 6  1008   248       94    69  3.60  4.81 Buckingham      34
 7  1011   195       92    41  4.80  4.84 Buckingham      30
 8  1015   227       75    44  5.20  3.94 Buckingham      37
 9  1016   177       87    49  3.60  4.84 Buckingham      45
10  1022   263       89    40  6.60  5.78 Buckingham      55
```

# Writing data

| Function | Writes |
|---|---|
| write_csv() | Comma separated values |
| write_excel_csv() | CSV that you plan to open in Excel |
| write_delim() | General delimited files |
| write_file() | A single string, written as is |
| write_lines() | A vector of strings, one string per line |
| write_tsv() | Tab delimited values |
| write_rds() | A data type used by R to save objects |
| write_xpt() | SAS transport format, .xpt |
| write_sas() | SAS .sas7bdat files (experimental) |

| Function | Writes |
|---|---|
| `write_sav()` | SPSS `.sav` files |
| `write_stata()` | Stata `.dta` files |

# Writing data

```
1  write_csv(diabetes, file = "diabetes-clean.csv")
```

# *Your Turn 4*

R has a few data file types, such as RDS and .Rdata. Save **diabetes** as **"diabetes.Rds"**.

...

```
1  write_rds(diabetes, "diabetes.Rds")
```