



THE UNIVERSITY *of* EDINBURGH

# Tidy(ing) data in R

---

Bram van Bunnik

Core Scientist  
Roslin Institute  
University of Edinburgh

# Tidy data - what is it & why do I need it?

- Huge amount of effort is spent cleaning data, estimates of up to 80% of analysis is spent on cleaning and preparing data
- Part of data cleaning is data tidying.
- Tidy datasets are easy to manipulate, model and visualise
- Tidy datasets have a specific structure:
- One table per dataset!!!



Illustrations from the Openscapes blog "Tidy Data for reproducibility, efficiency, and collaboration" by [Julia Lowndes](#) and [Allison Horst](#)

# Tidy data

## Tidy Data

*Tidy data* is a standard way of mapping the meaning of a dataset to its structure

Hadley Wickham

- Tidy datasets have a specific structure:
  - Each column is a variable
  - Each row is an observation
  - Each cell is a single value

each column a variable

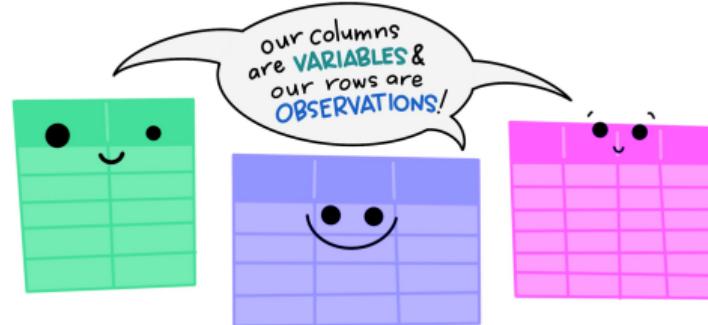
id	name	color
1	floof	gray
2	max	black
3	cat	orange
4	donut	gray
5	merlin	black
6	panda	calico

each row an observation



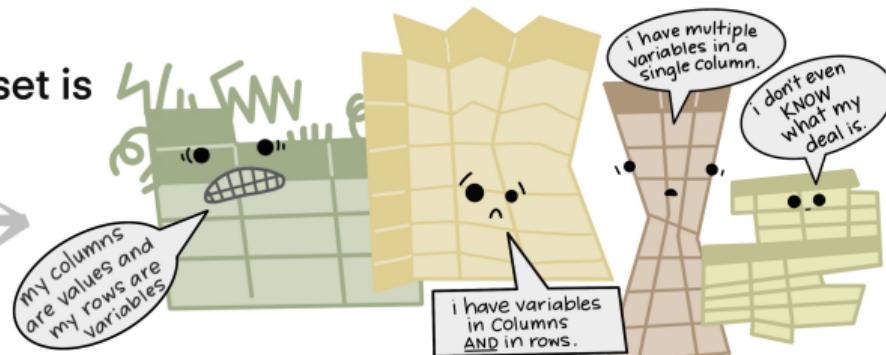
# All tidy datasets are alike...

The standard structure of  
tidy data means that  
“tidy datasets are all alike...”



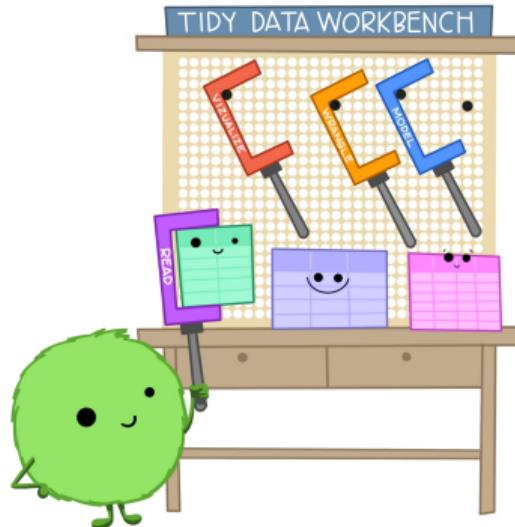
“...but every messy dataset is  
messy in its own way.”

-HADLEY WICKHAM



# Tidy datasets are easier to work with...

When working with tidy data,  
we can use the **same tools** in  
**similar ways** for different datasets...

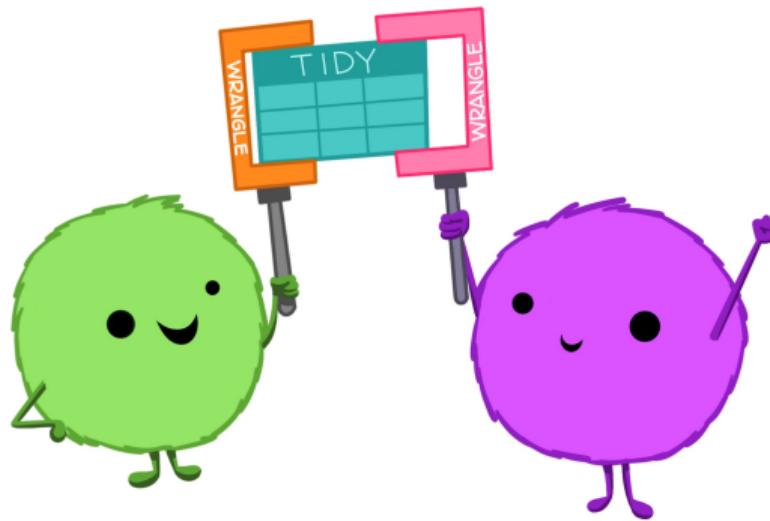


...but working with untidy data often means  
reinventing the wheel with **one-time**  
**approaches** that are **hard to iterate or reuse**.

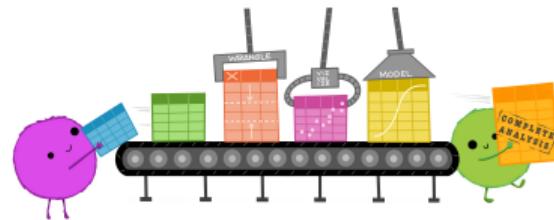


# Tidy datasets make sharing and collaborating easier...

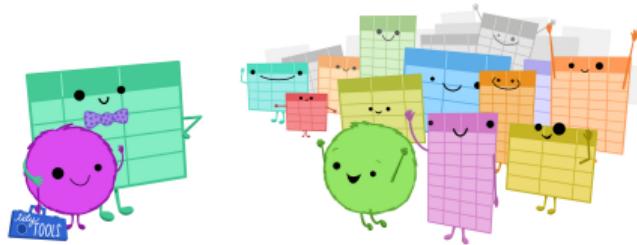
We all have the same tools we (can) use...



Tidy datasets are easier for automation and iteration

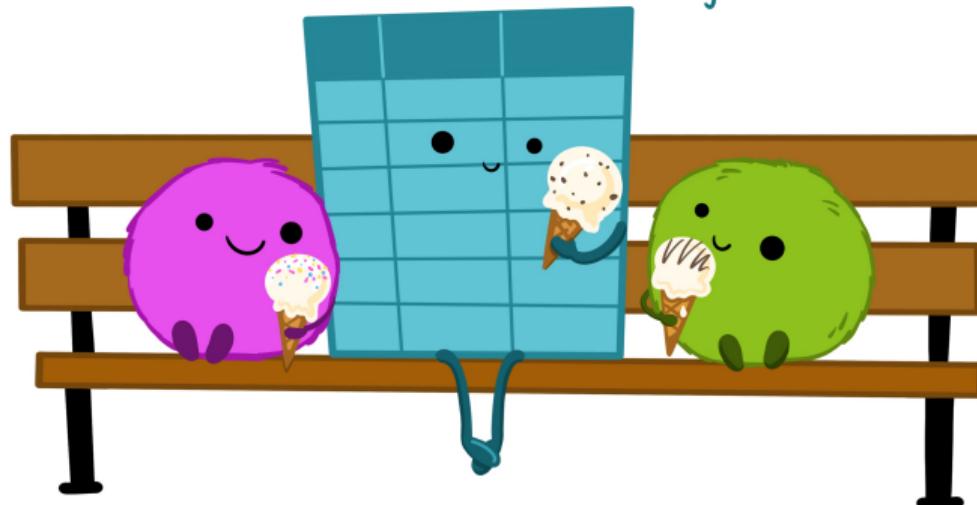


And it makes other (tidy) datasets seem more welcoming!



So...

make friends with tidy data.



## Some examples...

We have some data about an imaginary classroom in a format commonly seen in the wild. The table has three columns and four rows, and both rows and columns are labeled.

```
1 classroom
2 #> # A tibble: 4 × 4
3 #>   name    quiz1 quiz2 test1
4 #>   <chr>   <chr>  <chr> <chr>
5 #> 1 Billy    NA     D      C
6 #> 2 Suzy     F      NA     NA
7 #> 3 Lionel   B      C      B
8 #> 4 Jenny    A      A      B
```

There are many ways to structure the same underlying data. We could for example transpose the dataset and it would still show the same data:

```
1 classroom
2 #> # A tibble: 3 × 5
3 #>   assessment Billy Suzy Lionel Jenny
4 #>   <chr>       <chr> <chr> <chr> <chr>
5 #> 1 quiz1       NA     F      B      A
6 #> 2 quiz2       D      NA     C      A
7 #> 3 test1       C      NA     B      B
```

The data is the same, but the layout is different. Our vocabulary of rows and columns is simply not rich enough to describe why the two tables represent the same data. In addition to appearance, we need a way to describe the underlying semantics, or meaning, of the values displayed in the table.



- A dataset is a collection of **values**, usually either numbers (if quantitative) or strings (if qualitative).
- **Values** are organised in two ways:  
Every value belongs to a **variable** and an **observation**.
- A **variable** contains all values that measure the same underlying **attribute** (like height, temperature, duration) across units.
- An **observation** contains all values measured on the same **unit** (like a person, or a day, or a race) across attributes.

# Our example made tidy...

```
1 library(tidyverse)
2
3 classroom2 <- classroom |>
4   pivot_longer(quiz1:test1,
5     names_to = "assessment",
6     values_to = "grade") |>
7   arrange(name, assessment)
8
9 classroom2
10 #> # A tibble: 12 × 3
11 #>   name    assessment grade
12 #>   <chr>   <chr>      <chr>
13 #> 1 Billy   quiz1      NA
14 #> 2 Billy   quiz2      D
15 #> 3 Billy   test1      C
16 #> 4 Jenny   quiz1      A
17 #> 5 Jenny   quiz2      A
18 #> 6 Jenny   test1      B
19 #> 7 Lionel  quiz1      B
20 #> 8 Lionel  quiz2      C
21 #> 9 Lionel  test1      B
22 #> 10 Suzy   quiz1      F
23 #> # i 2 more rows
```

This makes the **values**, **variables**, and **observations** more clear.

The dataset contains 36 **values** representing 3 **variables** and 12 **observations**.

The **variables** are:

- "name", four possible values (Billy, Suzy, Lionel, and Jenny).
- "assessment", three possible values (quiz1, quiz2, and test1).
- "grade", five or six values (A, B, C, D, F, NA).



# Messy mistakes

---

Five most common "messy mistakes"

- Column headers are values, not variable names.
- Multiple variables are stored in one column.
- Variables are stored in both rows and columns.
- Multiple types of observational units are stored in the same table.
- A single observational unit is stored in multiple tables.

Surprisingly, most messy datasets, including types of messiness not explicitly described above, can be tidied with a small set of tools: pivoting (longer and wider) and separating.



# Column headers being values...

```
1 > chick_wide
2 # A tibble: 50 x 14
3   Chick Diet     T0     T2     T4     T6     T8     T10    T12    T14    T16    T18    T20    T21
4   <ord> <fct> <dbl> <dbl>
5   1 1     1     42     51     59     64     76     93    106    125    149    171    199    205
6   2 2     1     40     49     58     72     84    103    122    138    162    187    209    215
7   3 3     1     43     39     55     67     84     99    115    138    163    187    198    202
8   4 4     1     42     49     56     67     74     87    102    108    136    154    160    157
9   5 5     1     41     42     48     60     79    106    141    164    197    199    220    223
10  6 6     1     41     49     59     74     97    124    141    148    155    160    160    157
11  7 7     1     41     49     57     71     89    112    146    174    218    250    288    305
12  8 8     1     42     50     61     71     84    93    110    116    126    134    125    NA
13  9 9     1     42     51     59     68     85    96    90    92    93    100    100    98
14 10 10    1     41     44     52     63     74    81    89    96    101    112    120    124
15 # i 40 more rows
16 # i Use 'print(n = ...)' to see more rows
```

Data shows the weight of chicks (**Chick**) on different timepoints (**T0...T21**) for different diets (**Diet**).

In this dataset the columns starting with **T** are the timepoints at which a measurement was taken.

This data set is not tidy as each row has multiple observations...



# How to tidy?

Tidying this dataset is relatively straightforward:

```
1 chick_wide |>
2   pivot_longer(cols = -c(Chick, Diet),
3     names_to = "time",
4     values_to = "weight") |>
5     separate(time, into=c(NA, "time"), sep=1)
```

```
1 # A tibble: 600 x 4
2 Chick Diet    time weight
3 <ord> <fct>  <chr>  <dbl>
4 1 1      1        0     42
5 2 1      1        2     51
6 3 1      1        4     59
7 4 1      1        6     64
8 5 1      1        8     76
9 6 1      1       10     93
10 7 1      1       12    106
11 8 1      1       14    125
12 9 1      1       16    149
13 10 1     1       18    171
14 # i 590 more rows
15 # i Use 'print(n = ...)' to see more rows
```



# Variables are stored in both rows and columns

```
1 #> # A tibble: 5,769 × 22
2 #>   iso2    year    m04    m514    m014    m1524    m2534    m3544    m4554    m5564    m65     mu
3 #>   <chr> <int> <int>
4 #> 1 AD     1989    NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
5 #> 2 AD     1990    NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
6 #> 3 AD     1991    NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
7 #> 4 AD     1992    NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
8 #> 5 AD     1993    NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
9 #> 6 AD     1994    NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
10 #> 7 AD     1996    NA     NA     0      0      0      4      1      0      0      NA
11 #> 8 AD     1997    NA     NA     0      0      1      2      2      1      6      NA
12 #> 9 AD     1998    NA     NA     0      0      0      1      0      0      0      NA
13 #> 10 AD    1999   NA     NA     0      0      0      1      1      0      0      NA
14 #> # 5,759 more rows
15 #> # 10 more variables: f04 <int>, f514 <int>, f014 <int>, f1524 <int>,
16 #> #   f2534 <int>, f3544 <int>, f4554 <int>, f5564 <int>, f65 <int>,
17 #> #   fu <int>
```

- Dataset records the counts of confirmed tuberculosis cases by country, year, and demographic group.
- The demographic groups are broken down by **sex** (m, f) and **age** (0-14, 15-25, 25-34, 35-44, 45-54, 55-64, 65 and older, unknown).

How would you go about tidying this dataset?



# Multiple variables are stored in one column

Original dataset (kindly provided by Erhan Yalcindag):

```
1 A tibble: 386 x 83
2 ...1 ASV1 ASV2 ASV3 ASV4 ASV5 ASV6 ASV7 ASV8 ASV9 ASV10
3 <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
4 1 Querylength 233 233 292 281 281 281 283 292 292 292
5 2 1504-Jan20-Fec_1456_Calf 100 0 0 0 0 0 0 0 0 0
6 3 1504-Jan20-Fec_1868_Calf 0 100 0 0 0 0 0 0 0 0
7 4 1504-Jan20-Fec_216_Calf 0 0 0 0 0 0 0 0 0 0
8 5 1504-Jan20-Fec_396_Calf 0 0 100 0 0 0 0 0 0 0
9 6 1504-Nov19-Fec_1107_Calf 0 100 0 0 0 0 0 0 0 0
10 7 1504-Nov19-Fec_1172_Calf 0 0 0 100 0 0 0 0 0 0
11 8 1504-Nov19-Fec_1232_Calf 0 0 0 0 0 0 0 0 0 0
12 9 1504-Nov19-Fec_2434_Calf 0 0 0 0 100 0 0 0 0 0
13 10 1504-Nov19-Fec_278_Calf 100 0 0 0 0 0 0 0 0 0
14 # i 376 more rows
15 # i 72 more variables: ASV11 <dbl>, ASV12 <dbl>, ASV13 <dbl>,
```

Two (main) issues here:

- Multiple variables in one column.
- Multiple measurements in one row.

(Which other issue(s) can you spot?)



# Tidied data:

```
1 data |>
2   # remove first row
3   filter(!row_number() %in% c(1)) |>
4   rename("sample" = "...1") |>
5   # split sample into different columns
6   separate(col = sample,
7             into = c("sample", "date", "stage", "reads", "age")) |>
8   # re-order meaningfully
9   mutate(date = fct(date,
10         levels = c("Nov19", "Jan20", "May20", "July20")))) |>
11  mutate(stage = fct(stage)) |>
12  mutate(reads = as.numeric(reads)) |>
13  mutate(age = fct(age))
```

```
1 # A tibble: 385 x 87
2   sample date    stage reads age     ASV1  ASV2  ASV3  ASV4
3   <chr>  <fct>  <fct> <dbl> <fct> <dbl> <dbl> <dbl> <dbl>
4   1 1504 Jan20 Fec      1456 Calf     100     0     0     0
5   2 1504 Jan20 Fec      1868 Calf      0    100     0     0
6   3 1504 Jan20 Fec       216 Calf      0     0     0     0
7   4 1504 Jan20 Fec      396 Calf      0     0    100     0
8   5 1504 Nov19 Fec      1107 Calf      0    100     0     0
9   6 1504 Nov19 Fec      1172 Calf      0     0     0   100
10  7 1504 Nov19 Fec      1232 Calf      0     0     0     0
11  8 1504 Nov19 Fec      2434 Calf      0     0     0     0
12  9 1504 Nov19 Fec       278 Calf     100     0     0     0
13 10 1504 Nov19 Fec      7201 Calf      0     0     0     0
14 # i 375 more rows
15 # i 60 more variables: ASV5 <dbl>, ASV6 <dbl> ,...
```



# Tidy further by creating a "long" dataset

```
1 data |>
2   pivot_longer(cols=starts_with("ASV"),
3     names_to = "ASV", values_to = "identity") |>
4   mutate(ASV = fct(ASV))
```

```
1 # A tibble: 31,570 x 7
2 sample date   stage reads age   ASV      identity
3 <chr>   <fct>  <fct> <dbl> <fct> <fct>    <dbl>
4 1 1504   Jan20 Fec    1456 Calf   ASV1      100
5 2 1504   Jan20 Fec    1456 Calf   ASV2       0
6 3 1504   Jan20 Fec    1456 Calf   ASV3       0
7 4 1504   Jan20 Fec    1456 Calf   ASV4       0
8 5 1504   Jan20 Fec    1456 Calf   ASV5       0
9 6 1504   Jan20 Fec    1456 Calf   ASV6       0
10 7 1504   Jan20 Fec    1456 Calf  ASV7       0
11 8 1504   Jan20 Fec    1456 Calf  ASV8       0
12 9 1504   Jan20 Fec    1456 Calf  ASV9       0
13 10 1504  Jan20 Fec    1456 Calf ASV10      0
14 # i 31,560 more rows
15 # i Use 'print(n = ...)' to see more rows
```



- Excel... It has many functions that **don't play nicely** with data analysis.
- Colour coding/formatting cells to denote information... **This should be a variable**
- Merging cells... **Breaks the rules of variables / observations**
- Top row should always be **header names**... PLEASE, not blank rows for "aesthetics"
- Tidy data is rarely aesthetically pleasing to humans (and it is not intended to be!)
- Tidy data is specific to **your research question**

# Useful functions for tidying data

---

- `read_csv()`, `read_tsv()`, `read_*`(): read in text file in a specific format, returns a tibble
- `View()` or `glimpse()`: view the entire dataset
- `pivot_longer(data, cols, names_to = "name", values_to = "value", values_drop_na = FALSE)`: Lengthen data by collapsing several columns into two
- `pivot_wider()`: the opposite of `pivot_longer`
- `separate()`, `separate_wider_delim()`, `separate_wider_position()`: turns a single character column into multiple columns
- `separate_longer_*`(): turn into multiple rows
- `mutate()`: create new column(s) that are functions of existing variables. It can also modify or delete columns



## Useful functions for tidying data, cont.

---

- `select()`: select variables in a dataframe
- `filter()`: subset dataframe to only include rows that satisfy condition
- `bind_rows()`, `bind_cols()`: bind any number of df's by row / column
- `{inner, left, right, full}_join()`: Mutating joins add columns from y to x, matching observations based on the keys. There are four mutating joins: the inner join, and the three outer joins.
- `clean_names()` (from janitor package): Cleans names of a df. Resulting names are unique and consist only of the "\_" character, numbers, and letters
- Any others??



Various cheatsheets from R-Studio (and others):

<https://github.com/rstudio/cheatsheets/>

Some noteworthy ones:

- <https://github.com/rstudio/cheatsheets/blob/main/pngs/tidyr.png>
- <https://github.com/rstudio/cheatsheets/blob/main/pngs/rstudio-ide.png>
- <https://github.com/rstudio/cheatsheets/blob/main/pngs/gtsummary.png>
- <https://github.com/rstudio/cheatsheets/blob/main/pngs/factors.png>

And remember:

make friends with tidy data.

