

Homework 5

Carmen Canedo

2 October 2020

Exercise 2.1 Why does tidy data lend itself to vectorised operations?

Tidy data ensures that an observation is always correctly paired with the variables

Exercise 2.2 How could you tidy the SAT data from last week? Which of the data sets below are tidy? What's wrong with the non-tidy data sets?

After reading in the SAT data from the .csv, I placed my data into a tibble. Each variable is in a column, each observation has its own row, and each value has its own cell. To make the SAT data tidy, I would also make all variable names use snake case (all lowercase letters and underscore instead of spaces), and use title case for all the name of the high schools for consistency.

The only table that is tidy:

- Table 1

The following **are not** tidy:

- Table 2
 - `rate` contains two variables. To fix this we can separate them into `num_cases` and `total_population`, and if we still wanted to include the rate, we could use `mutate()` to divide the two and store the values in `rate`.
- Table 3
 - 2000 and 1999 belong to one variable `year`, but in this table, they are spread across two columns. To fix this, we can use `pivot_longer()` and assign column names to `year` and the values to a separate column, `num_cases`.
- Table 4
 - The observations (country names) in the rows are repeated, so we can use `pivot_wider()` to split `type` into `num_cases` and `total_population`

Exercise 2.3 Use `pivot_longer()` to tidy data frame

```
## # A tibble: 6 x 11
##   religion '$10k' '$10-20k' '$20-30k' '$30-40k' '$40-50k' '$50-75k' '$75-100k'
##   <chr>      <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>
## 1 Agnostic    27      34      60      81      76     137     122
## 2 Atheist     12      27      37      52      35      70      73
```

```
## 3 Buddhist      27      21      30      34      33      58      62
## 4 Catholic     418     617     732     670     638     1116     949
## 5 Don't k~      15      14      15      11      10      35      21
## 6 Evangel~     575     869     1064     982     881     1486     949
## # ... with 3 more variables: '$100-150k' <dbl>, '>150k' <dbl>, 'Don't
## #   know/refused' <dbl>
```

```
## # A tibble: 180 x 3
##   religion income      count
##   <chr>    <chr>    <dbl>
## 1 Agnostic <$10k      27
## 2 Agnostic $10-20k    34
## 3 Agnostic $20-30k    60
## 4 Agnostic $30-40k    81
## 5 Agnostic $40-50k    76
## 6 Agnostic $50-75k   137
## 7 Agnostic $75-100k  122
## 8 Agnostic $100-150k 109
## 9 Agnostic >150k     84
## 10 Agnostic Don't know/refused 96
## # ... with 170 more rows
```

Exercise 2.4 Tidy the data from blackboard

```
## # A tibble: 90 x 4
##   Monkey Treatment week  accuracy
##   <chr>    <chr>    <chr>    <dbl>
## 1 Spank   Control   Week2     95
## 2 Spank   Control   Week4     75
## 3 Spank   Control   Week8     80
## 4 Spank   Control   Week12    65
## 5 Spank   Control   Week16    70
## 6 Chim    Control   Week2     85
## 7 Chim    Control   Week4     75
## 8 Chim    Control   Week8     55
## 9 Chim    Control   Week12    75
## 10 Chim   Control   Week16    85
## # ... with 80 more rows
```

Exercise 2.5 Use pivot_wider() to tidy tidyr::fishencounters

```
## # A tibble: 19 x 12
##   fish Release I80_1 Lisbon Rstr Base_TD BCE BCW BCE2 BCW2 MAE MAW
##   <fct>    <int> <int>    <int> <int>    <int> <int> <int> <int> <int> <int> <int>
## 1 4842      1      1      1      1      1      1      1      1      1      1      1
## 2 4843      1      1      1      1      1      1      1      1      1      1      1
## 3 4844      1      1      1      1      1      1      1      1      1      1      1
## 4 4845      1      1      1      1      1      NA     NA     NA     NA     NA     NA
## 5 4847      1      1      1      NA     NA     NA     NA     NA     NA     NA     NA
## 6 4848      1      1      1      1      NA     NA     NA     NA     NA     NA     NA
## 7 4849      1      1      NA     NA     NA     NA     NA     NA     NA     NA     NA
## 8 4850      1      1      NA     1      1      1      1      NA     NA     NA     NA
```

```
## 9 4851      1      1      NA      NA      NA      NA      NA      NA      NA      NA      NA
## 10 4854     1      1      NA      NA      NA      NA      NA      NA      NA      NA      NA
## 11 4855     1      1      1      1      1      NA      NA      NA      NA      NA      NA
## 12 4857     1      1      1      1      1      1      1      1      1      NA      NA
## 13 4858     1      1      1      1      1      1      1      1      1      1      1
## 14 4859     1      1      1      1      1      NA      NA      NA      NA      NA      NA
## 15 4861     1      1      1      1      1      1      1      1      1      1      1
## 16 4862     1      1      1      1      1      1      1      1      1      NA      NA
## 17 4863     1      1      NA      NA      NA      NA      NA      NA      NA      NA      NA
## 18 4864     1      1      NA      NA      NA      NA      NA      NA      NA      NA      NA
## 19 4865     1      1      1      NA      NA      NA      NA      NA      NA      NA      NA
```

Exercise 2.6 Tidy flowers1 data set

```
## # A tibble: 24 x 4
##   Time replication Flowers Intensity
##   <dbl> <chr>      <chr>    <chr>
## 1     1 " 1"        " 62,3" 150,0
## 2     1 " 2"        " 77,4" 150,0
## 3     1 " 3"        " 55,3" 300,0
## 4     1 " 4"        " 54,2" 300,0
## 5     1 " 5"        " 49,6" 450,0
## 6     1 " 6"        " 61,9" 450,0
## 7     1 " 7"        " 39,4" 600,0
## 8     1 " 8"        " 45,7" 600,0
## 9     1 " 9"        " 31,3" 750,0
## 10    1 "10"       " 44,9" 750,0
## # ... with 14 more rows
```

Exercise 2.7 Use separate to tidy the flowers2 data set

```
## # A tibble: 24 x 3
##   flowers intensity Time
##   <dbl>      <int> <dbl>
## 1    62.3      150     1
## 2    77.4      150     1
## 3    55.3      300     1
## 4    54.2      300     1
## 5    49.6      450     1
## 6    61.9      450     1
## 7    39.4      600     1
## 8    45.7      600     1
## 9    31.3      750     1
## 10   44.9      750     1
## # ... with 14 more rows
```

Exercise 2.8 Read the help file for unite and correct the code above to get rid of underscore in year column

```
## # A tibble: 6 x 3
##   country    year rate
```

```
##   <chr>      <chr> <chr>
## 1 Afghanistan 1999  745/19987071
## 2 Afghanistan 2000  2666/20595360
## 3 Brazil      1999  37737/172006362
## 4 Brazil      2000  80488/174504898
## 5 China       1999  212258/1272915272
## 6 China       2000  213766/1280428583
```

Exercise 2.9 Turn implicit missing values explicit in the data frame

```
## # A tibble: 5 x 3
##   treatment gender return
##   <chr>      <fct>   <dbl>
## 1 a        M        1.5
## 2 b        F        0.75
## 3 a        F        0.5
## 4 c        M        1.8
## 5 b        M        NA
```

```
## # A tibble: 15 x 3
##   treatment gender return
##   <chr>      <fct>   <dbl>
## 1 a        M        1.5
## 2 a        M        1.8
## 3 a        M        NA
## 4 a        F        0.5
## 5 a        F        0.75
## 6 b        M        1.5
## 7 b        M        1.8
## 8 b        M        NA
## 9 b        F        0.5
## 10 b       F        0.75
## 11 c       M        1.5
## 12 c       M        1.8
## 13 c       M        NA
## 14 c       F        0.5
## 15 c       F        0.75
```

Exercise 2.10 Tidy the tidyr::billboard data set

- 1: Gather up all the week entries into a row for each week for each song where there is an entry
- 2: Convert the week variable to a number and figure out the date corresponding to each week on the chart
- 3: Sort the data by artist, track, and week

```
## # A tibble: 5,307 x 5
##   artist track week position date
##   <chr> <chr>   <dbl>   <dbl> <date>
## 1 2 Pac  Baby Don't Cry (Keep... 1      87 2000-02-26
```

##	2	2	Pac	Baby Don't Cry (Keep...	2	82	2000-03-04
##	3	2	Pac	Baby Don't Cry (Keep...	3	72	2000-03-11
##	4	2	Pac	Baby Don't Cry (Keep...	4	77	2000-03-18
##	5	2	Pac	Baby Don't Cry (Keep...	5	87	2000-03-25
##	6	2	Pac	Baby Don't Cry (Keep...	6	94	2000-04-01
##	7	2	Pac	Baby Don't Cry (Keep...	7	99	2000-04-08
##	8	2	Ge+her	The Hardest Part Of ...	1	91	2000-09-02
##	9	2	Ge+her	The Hardest Part Of ...	2	87	2000-09-09
##	10	2	Ge+her	The Hardest Part Of ...	3	92	2000-09-16
##	#			... with 5,297 more rows			