

Tidy Data

Byteflow Dynamics
9/24/2017

Why tidy data?

- Easier to use the tools in the tidyverse package
- You will spend less time munging data
- And more time on analytics

tidyr

We will use tidyr (part of tidyverse) to organize messy data

```
library(tidyverse)

## Loading tidyverse: ggplot2
## Loading tidyverse: tibble
## Loading tidyverse: tidyr
## Loading tidyverse: readr
## Loading tidyverse: purrr
## Loading tidyverse: dplyr

## Warning: package 'tidyr' was built under R version 3.4.1
## Warning: package 'purrr' was built under R version 3.4.1
## Conflicts with tidy packages -----
## filter(): dplyr, stats
## lag():    dplyr, stats
```

What is tidy data?

1. Each variable must have its own column.
2. Each observation must have its own row.
3. Each value must have its own cell.

Example: Load tables 1, 2, 3, 4a, and 4b and observe differences.

```
table1

## # A tibble: 6 x 4
##   country year cases population
##   <chr> <int> <int>      <int>
## 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
```

table2

```
## # A tibble: 12 x 4
##   country year   type    count
##   <chr> <int>   <chr>   <int>
## 1 Afghanistan 1999   cases     745
## 2 Afghanistan 1999 population 19987071
## 3 Afghanistan 2000   cases     2666
## 4 Afghanistan 2000 population 20595360
## 5      Brazil 1999   cases    37737
## 6      Brazil 1999 population 172006362
## 7      Brazil 2000   cases     80488
## 8      Brazil 2000 population 174504898
## 9       China 1999   cases    212258
## 10      China 1999 population 1272915272
## 11      China 2000   cases    213766
## 12      China 2000 population 1280428583
```

table3

```
## # A tibble: 6 x 3
##   country year      rate
## *   <chr> <int>   <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
```

table4a

```
## # A tibble: 3 x 3
##   country `1999` `2000`
## *   <chr> <int> <int>
## 1 Afghanistan    745    2666
## 2      Brazil 37737    80488
## 3       China 212258    213766
```

table4b

```
## # A tibble: 3 x 3
##   country `1999` `2000`
## *   <chr>   <int>   <int>
## 1 Afghanistan 19987071 20595360
## 2      Brazil 172006362 174504898
## 3       China 1272915272 1280428583
```

Question 1: Which one(s) is/are tidy?

Question 2: Are our grade datasets tidy?

How to make data tidy?

Data can be messy in 3 ways. For each problem, there is a tidyr function.

1. One variable is spread across multiple columns (tables4a, 4b)

- Use `gather()`
2. One observation is scattered across multiple rows (table2)
 - Use `spread()`
 3. A single cell contains more than one values (table3)
 - Use `separate()`

How do they work?

1. First, pick a data frame
2. Then apply a function, describing what to do with the data frame
3. The result is a new data frame

Pipe

The pipe operator `%>%` is used to make codes more readable

gather()

If a dataset has a variable spread accross multiple columns, we need to gather these columns.

```
table4a
```

```
## # A tibble: 3 x 3
##   country `1999` `2000`
## *   <chr>   <int> <int>
## 1 Afghanistan    745   2666
## 2      Brazil  37737  80488
## 3        China 212258 213766
```

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

```
## # A tibble: 6 x 3
##   country year cases
##   <chr> <chr> <int>
## 1 Afghanistan 1999    745
## 2      Brazil 1999  37737
## 3        China 1999 212258
## 4 Afghanistan 2000   2666
## 5      Brazil 2000  80488
## 6        China 2000 213766
```

These produce the same results

```
table4a %>%
  gather(2:3, key = "year", value = "cases")
```

```
## # A tibble: 6 x 3
##   country year cases
##   <chr> <chr> <int>
## 1 Afghanistan 1999    745
## 2      Brazil 1999  37737
## 3        China 1999 212258
```

```
## 4 Afghanistan 2000 2666
## 5      Brazil 2000 80488
## 6      China 2000 213766
```

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

```
## # A tibble: 6 x 3
##   country year cases
##   <chr> <chr> <int>
## 1 Afghanistan 1999 745
## 2      Brazil 1999 37737
## 3      China 1999 212258
## 4 Afghanistan 2000 2666
## 5      Brazil 2000 80488
## 6      China 2000 213766
```

```
table4a %>%
  gather(-1, key = "year", value = "cases")
```

```
## # A tibble: 6 x 3
##   country year cases
##   <chr> <chr> <int>
## 1 Afghanistan 1999 745
## 2      Brazil 1999 37737
## 3      China 1999 212258
## 4 Afghanistan 2000 2666
## 5      Brazil 2000 80488
## 6      China 2000 213766
```

Exercise: Use `gather()` to tidy up table 4b. Use “population” as the value.

```
table4b %>%
  gather(`1999`, `2000`, key = "year", value = "population")
```

```
## # A tibble: 6 x 3
##   country year population
##   <chr> <chr>      <int>
## 1 Afghanistan 1999 19987071
## 2      Brazil 1999 172006362
## 3      China 1999 1272915272
## 4 Afghanistan 2000 20595360
## 5      Brazil 2000 174504898
## 6      China 2000 1280428583
```

spread()

When an observation is scattered across multiple rows, we need to spread them.

```
table2
```

```
## # A tibble: 12 x 4
##   country year      type count
##   <chr> <int>    <chr> <int>
## 1 Afghanistan 1999 cases      745
## 2 Afghanistan 1999 population 19987071
```

```
## 3 Afghanistan 2000 cases 2666
## 4 Afghanistan 2000 population 20595360
## 5 Brazil 1999 cases 37737
## 6 Brazil 1999 population 172006362
## 7 Brazil 2000 cases 80488
## 8 Brazil 2000 population 174504898
## 9 China 1999 cases 212258
## 10 China 1999 population 1272915272
## 11 China 2000 cases 213766
## 12 China 2000 population 1280428583
```

```
table2 %>%
  spread(key = type, value = count)
```

```
## # A tibble: 6 x 4
##   country year cases population
## *   <chr> <int> <int>    <int>
## 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
```

separate()

If there are multiple values in a cell, we separate them.

```
table3
```

```
## # A tibble: 6 x 3
##   country year      rate
## *   <chr> <int>    <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
```

```
table3 %>%
  separate(rate, into = c("cases", "population"), sep = "/")
```

```
## # A tibble: 6 x 4
##   country year cases population
## *   <chr> <int> <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: Tidy the lab data

Hint: use key = “type” and value = “score”

```
tidylab <- read.csv("lab_grades.csv", stringsAsFactors = FALSE) %>%  
  gather(Lab.1:Lab.14, key = "type", value = "score")  
  
write.csv(tidylab, "tidylab.csv")
```