Data munging

Abhijit Dasgupta

March 25-27, 2019

Tidy data

Make the computer happy

Just as we want code to make humans happy, we want data to make computers happy

Tidy data is a principle promoted by Dr. Hadley Wickham to make computers happy

The properties of a tidy dataset are:

- 1. Each variable forms a column
- 2. Each observation forms a row
- 3. Each type of observational unit forms a table.

This forms a standardized way to structure a dataset, and so makes it easy for the analyst to develop standard pipelines.

A tidy dataset is tidy in one way, but a messy dataset can be messy in many ways

Hadley Wickham

Messy data

A dataset can be messy in many many ways. Many of the more common issues are listed below:

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

Sometimes the messier format is better for data entry, but bad for data analyses.

Variables in column names

```
library(tidyverse)
pew <- import('Data/FSI/pew.csv')</pre>
head(pew, 4)
  religion <$10k $10-20k $20-30k $30-40k $40-50k $50-75k $75-100k
                                                 76
1 Agnostic
                                                        137
                                                                  122
 Atheist
                                                 35
                                                                   73
3 Buddhist
                                                33
                                                                   62
4 Catholic
                                                638
                                                       1116
                                                                  949
  $100-150k >150k Don't know/refused
        109
               84
                                    76
               53
                                    54
        792
              633
                                  1489
```

- This dataset has actual data in the column headers, rather than variable names.
- We should ideally have 3 columns in this dataset: religion, income and frequency.
- We can achieve this using a function called gather which takes a wide dataset and makes it tall.

- Gather all the columns into two columns, income and frequency, by stacking the columns
- Don't include the variable religion

```
pew %>%
  gather(income, frequency, -religion) %>%
  as_tibble()
```

```
# A tibble: 180 x 3
                            income frequency
   religion
   <chr>
                            <chr>
                                       <int>
 1 Agnostic
                            <$10k
                                          27
 2 Atheist
                                          12
                            <$10k
 3 Buddhist
                            <$10k
                                          27
 4 Catholic
                            <$10k
                                         418
 5 Don't know/refused
                            <$10k
                                          15
                                         575
 6 Evangelical Prot
                            <$10k
 7 Hindu
                            <$10k
 8 Historically Black Prot <$10k
                                         228
 9 Jehovah's Witness
                            <$10k
                                          20
10 Jewish
                            <$10k
                                          19
# ... with 170 more rows
```

PS 312, March 2019

head(pew)

| | | religion | <\$10k | \$10-20k | \$20-30k | \$30-40k | \$4 |
|---|-----------|------------|--------|----------|-----------|----------|-----|
| 1 | | Agnostic | 27 | 34 | 60 | 81 | |
| 2 | | Atheist | 12 | 27 | 37 | 52 | |
| 3 | | Buddhist | 27 | 21 | 30 | 34 | |
| 4 | | Catholic | 418 | 617 | 732 | 670 | |
| 5 | Don't kno | ow/refused | 15 | 14 | 15 | 11 | |
| 6 | Evange] | lical Prot | 575 | 869 | 1064 | 982 | |
| | \$75-100k | \$100-150k | >150k | Don't kr | now/refus | sed | |
| 1 | 122 | 109 | 84 | | | 96 | |
| 2 | 73 | 59 | 74 | | | 76 | |
| 3 | 62 | 39 | 53 | | | 54 | |
| 4 | 949 | 792 | 633 | | 14 | 489 | |
| 5 | 21 | 17 | 18 | | • | 116 | |
| 6 | 949 | 723 | 414 | | 15 | 529 | |

pew %>% gather(income, frequency, -religion) %>% head(20)

| | religion | | frequency |
|-----|-------------------------|--------|-----------|
| 1 | Agnostic | <\$10k | 27 |
| 2 | Atheist | <\$10k | 12 |
| 3 | Buddhist | <\$10k | 27 |
| 4 | Catholic | <\$10k | 418 |
| 5 | Don't know/refused | <\$10k | 15 |
| 6 | Evangelical Prot | <\$10k | 575 |
| 7 | Hindu | <\$10k | 1 |
| 8 F | Historically Black Prot | <\$10k | 228 |
| 9 | Jehovah's Witness | <\$10k | 20 |
| 10 | Jewish | <\$10k | 19 |
| 11 | Mainline Prot | <\$10k | 289 |
| 12 | Mormon | <\$10k | 29 |
| 13 | Muslim | <\$10k | 6 |
| 14 | Orthodox | <\$10k | 13 |
| 15 | Other Christian | <\$10k | 9 |
| 16 | Other Faiths | <\$10k | 20 |
| 17 | Other World Religions | <\$10k | 5 |
| 18 | Unaffiliated | | 217 |
| 19 | Agnostic | | 34 |
| 20 | Atheist | | 27 |
| | | | |

Multiple variables in column names

```
tb <- import('Data/FSI/tb.csv') %>% as_tibble()
head(tb)
```

```
# A tibble: 6 x 22
 iso2
      vear
            m04
                m514 m014 m1524 m2534 m3544 m4554 m5564
 1 AD
      1989
             NA
                           NA
                                NA
                                                        NA
2 AD
     1990
             NA
                           NA
                                                        NA
3 AD
     1991
             NA NA
                           NA
                                                        NA
                           NA NA
4 AD
     1992
           NA NA
                                                        NA
                           NA
5 AD
     1993
             NA
                                                        NA
6 AD
      1994
                           NA
                                                        NA
# ... with 10 more variables: f04 <int>, f514 <int>, f014 <int>,
   f1524 <int>, f2534 <int>, f3544 <int>, f4554 <int>, f5564 <int>,
   f65 <int>, fu <int>
```

column headers include both sex and age

```
tb %>%
  gather(sex_age, n, -iso2, -year, -fu)
```

```
# A tibble: 109,611 x 5
  iso2 year fu sex_age
                              n
  <chr> <int> <int> <chr> <int>
1 AD
         1989
                 NA m04
                              NA
2 AD
         1990
                NA m04
                              NA
3 AD
         1991
                 NA m04
                              NA
4 AD
         1992
                 NA m04
                              NA
5 AD
         1993
                 NA m04
                              NA
6 AD
         1994
                 NA m04
                              NA
7 AD
         1996
                 NA m04
                              NA
8 AD
         1997
                 NA m04
                              NA
9 AD
         1998
                 NA m04
                              NA
10 AD
         1999
                 NA m04
                              NA
# ... with 109,601 more rows
```

```
tb %>%
  gather(sex_age, n, -iso2, -year, -fu, na.rm=T)
```

```
# A tibble: 35,478 x 5
  iso2 year fu sex_age
                              n
  <chr> <int> <int> <chr> <int>
1 AD
         2005
                 0 m04
2 AD
         2006
                0 m04
3 AD
         2008
                0 m04
4 AE
         2006
                NA m04
5 AE
         2007
                NA m04
6 AE
         2008
                0 m04
7 AG
         2007
                NA m04
8 AL
         2005
                0 m04
9 AL
         2006
                 0 m04
10 AL
         2007
                 0 m04
                               0
# ... with 35,468 more rows
```

```
tb %>%
  gather(sex_age, n, -iso2, -year, -fu, na.rm=T) %>%
separate(sex_age, c('sex', 'age'), sep=1) # by position
# A tibble: 35,478 x 6
   iso2 year fu sex
                              age
                                          n
   <chr> <int> <int> <chr> <chr> <int>
 1 AD
           2005
                     0 m
                              04
 2 AD
           2006
                     0 m
                              04
 3 AD
           2008
                     0 m
                              04
 4 AE
           2006
                    NA m
                              04
 5 AE
           2007
                    NA m
                              04
 6 AE
           2008
                    0 m
                              04
 7 AG
           2007
                    NA m
                              04
 8 AL
           2005
                     0 m
                              04
 9 AL
           2006
                     0 m
                              04
10 AL
           2007
                     0 m
                              04
                                          0
# ... with 35,468 more rows
```

This still needs to be cleaned

Variables stored in rows and columns

```
weather <- import('Data/FSI/weather.csv') %>% as_tibble()
weather
```

```
# A tibble: 22 x 35
         id
                                 vear month element
                                                                                                           d1
                                                                                                                                d2
                                                                                                                                                    d3
                                                                                                                                                                         d4
                                                                                                                                                                                             d5
                                                                                                                                                                                                                                     d7
          <chr> <int> <int> <chr> <dbl> <
   1 MX17... 2010
                                                           1 tmax
                                                                                                           NA NA
                                                                                                                                             NA
                                                                                                                                                                        NA NA
   2 MX17... 2010
                                                          1 tmin
                                                                                                           NA NA
                                                                                                                                                                        NA NA
   3 MX17... 2010
                                                             2 tmax
                                                                                                           NA 27.3
                                                                                                                                             24.1
   4 MX17... 2010
                                                             2 tmin
                                                                                                           NA 14.4
                                                                                                                                             14.4
                                                                                                                                                                        NA NA
   5 MX17... 2010
                                                                                                                                                                        NA 32.1
                                                            3 tmax
   6 MX17... 2010
                                                             3 tmin
                                                                                                         NA NA
                                                                                                                                                                        NA 14.2
                                                                                                                                                                                                                                     NA
   7 MX17... 2010
                                                                                                         NA NA
                                                            4 tmax
                                                                                                                                                                        NA NA
                                                                                                                                                                                                                                     NA
                                                                                           NA NA
   8 MX17... 2010
                                                          4 tmin
                                                                                                                                             NA
                                                                                                                                                                        NA NA
                                                                                                                                                                                                                                     NA
   9 MX17... 2010
                                                             5 tmax
                                                                                                         NA NA
                                                                                                                                             NA
                                                                                                                                                                        NA NA
10 MX17... 2010
                                                               5 tmin
                                                                                                           NA NA
                                                                                                                                             NA
                                                                                                                                                                        NA NA
# ... with 12 more rows, and 24 more variables: d8 <dbl>, d9 <lgl>,
             d10 <dbl>, d11 <dbl>, d12 <lgl>, d13 <dbl>, d14 <dbl>, d15 <dbl>,
           d16 <dbl>, d17 <dbl>, d18 <lgl>, d19 <lgl>, d20 <lgl>, d21 <lgl>,
          d22 <1g1>, d23 <db1>, d24 <1g1>, d25 <db1>, d26 <db1>, d27 <db1>,
             d28 <db1>, d29 <db1>, d30 <db1>, d31 <db1>
```

```
weather %>%
  gather(day, temp, d1:d31)
```

```
# A tibble: 682 x 6
  id
            year month element day
                                        temp
   <chr>
           <int> <int> <chr>
                                <chr> <dbl>
 1 MX17004
            2010
                       tmax
                                d1
                                          NA
 2 MX17004
            2010
                      1 tmin
                                d1
                                          NA
 3 MX17004
            2010
                      2 tmax
                                d1
                                          NA
 4 MX17004
            2010
                      2 tmin
                                d1
                                          NA
 5 MX17004
            2010
                                d1
                                          NA
                      3 tmax
 6 MX17004
            2010
                      3 tmin
                                d1
                                          NA
7 MX17004
            2010
                      4 tmax
                                d1
                                          NA
 8 MX17004
            2010
                      4 tmin
                                d1
                                          NA
 9 MX17004
            2010
                      5 tmax
                                d1
                                          NA
10 MX17004
            2010
                      5 tmin
                                d1
                                          NA
# ... with 672 more rows
```

d1:d31 denotes all the variables physically between d1 and d31 in the dataset

See what happens when you type 1:10 in the console

```
weather %>%
  gather(date, temp, d1:d31) %>%
  spread(element, temp)
```

```
# A tibble: 341 x 6
  id
           year month date
                             tmax tmin
   <chr>
          <int> <int> <chr> <dbl> <dbl>
 1 MX17004
                     1 d1
           2010
                                NA
                                     NA
2 MX17004
           2010
                    1 d10
                                NA
                                     NA
3 MX17004
           2010
                    1 d11
                                NA
                                     NA
4 MX17004
           2010
                    1 d12
                                NA
                                     NA
 5 MX17004
           2010
                    1 d13
                                NA
                                     NA
 6 MX17004
           2010
                     1 d14
                                NA
                                     NA
 7 MX17004
           2010
                     1 d15
                                     NA
                                NA
8 MX17004
           2010
                    1 d16
                                NA
                                     NA
9 MX17004
                     1 d17
           2010
                                NA
                                     NA
10 MX17004
           2010
                     1 d18
                                NA
                                     NA
# ... with 331 more rows
```

Data cleaning

```
weather %>%
  gather(date, temp, d1:d31) %>%
  spread(element, temp)
# A tibble: 341 x 6
  id
           year month date
                              tmax tmin
   <chr>
           <int> <int> <chr> <dbl> <dbl>
 1 MX17004
                     1 d1
           2010
                                NA
                                      NA
 2 MX17004
           2010
                     1 d10
                                NA
                                      NA
 3 MX17004
           2010
                     1 d11
                                NA
                                      NA
 4 MX17004
           2010
                     1 d12
                                NA
                                      NA
 5 MX17004
           2010
                     1 d13
                                      NA
                                NA
 6 MX17004
           2010
                     1 d14
                                NA
                                      NA
 7 MX17004
                     1 d15
                                      NA
           2010
                                NA
 8 MX17004
           2010
                     1 d16
                                NA
                                      NA
9 MX17004
           2010
                     1 d17
                                NA
                                      NA
10 MX17004
           2010
                     1 d18
                                      NA
                                NA
# ... with 331 more rows
```

• date column in alphabetical rather than numerical order

```
weather %>%
  gather(date, temp, d1:d31) %>%
  spread(element, temp) %>%
  mutate(date = parse_number(date))
```

```
# A tibble: 341 x 6
  id
           year month date tmax tmin
          <int> <int> <dbl> <dbl> <dbl>
  <chr>
 1 MX17004
          2010
                               NA
                                    NA
 2 MX17004 2010
                         10
                               NA
                                    NA
 3 MX17004
          2010
                         11
                               NA
                                    NA
 4 MX17004
          2010
                         12
                                    NA
                               NA
 5 MX17004
          2010
                         13
                               NA
                                    NA
 6 MX17004
          2010
                         14
                                    NA
                               NA
 7 MX17004
          2010
                         15
                               NA
                                    NA
 8 MX17004
          2010
                         16
                               NA
                                    NA
9 MX17004
          2010
                         17
                                    NA
                               NA
10 MX17004
          2010
                         18
                               NA
                                    NA
# ... with 331 more rows
```

```
weather %>%
  gather(date, temp, d1:d31) %>%
  spread(element, temp) %>%
  mutate(date = parse_number(date)) %>%
  arrange(date)
```

```
# A tibble: 341 x 6
  id
           year month date tmax tmin
          <int> <int> <dbl> <dbl> <dbl>
   <chr>
 1 MX17004
           2010
                               NA
                                     NA
 2 MX17004 2010
                               NA
                                     NA
 3 MX17004 2010
                                     NA
                               NA
 4 MX17004
          2010
                               NA
                                     NA
 5 MX17004
                                     NA
          2010
                               NA
 6 MX17004
          2010
                               NA
                                     NA
 7 MX17004
          2010
                               NA
                                     NA
 8 MX17004
          2010
                               NA
                                     NA
 9 MX17004
          2010
                               NA
                                     NA
10 MX17004
           2010
                               NA
                                     NA
# ... with 331 more rows
```

Not quite! We'd like dates ordered within months

```
weather %>%
  gather(date, temp, d1:d31) %>%
  spread(element, temp) %>%
 mutate(date = parse_number(date)) %>%
 arrange(month, date)
# A tibble: 341 x 6
           vear month date tmax tmin
  id
   <chr> <int> <int> <dbl> <dbl> <dbl>
 1 MX17004 2010
                               NA
                                     NA
 2 MX17004 2010
                               NA
                                     NA
 3 MX17004 2010
                               NA
                                     NA
 4 MX17004 2010
                                     NA
                               NA
5 MX17004 2010
                                     NA
                               NA
 6 MX17004 2010
                               NA
                                     NA
 7 MX17004 2010
                               NA
                                     NA
8 MX17004 2010
                                     NA
                               NA
9 MX17004 2010
                               NA
                                     NA
10 MX17004 2010
                                     NA
# ... with 331 more rows
```

Good. Now to save it

```
weather2 <- weather %>%
  gather(date, temp, d1:d31) %>%
  spread(element, temp) %>%
  mutate(date = parse_number(date)) %>%
  arrange(month, date)
```

Exercise

The file Data/FSI/mbta.xlsx contains monthly data on number of commuter trips by different modalities on the MBTA system n Boston.

- It is in a messy format.
- It also has an additional quirk in that it has a title on the first line that isn't even data. You can avoid loading that in by using the option skip=1 (i.e. skip the first line) when you import.

Work through this process to clean this dataset into tidy form. I'll also note that you can "minus" columns by position as well as name, so gather(date, avg_trips, -1, -mode) is valid to not involve the first column and the mode column.

```
mbta <- import('Data/FSI/mbta.xlsx', skip = 1) %>% as_tibble()
mbta2 <- mbta %>%
   gather(date, avg_trips, -1, -mode) %>%
   separate(date, c("year", "month"), sep = '-')
mbta2
```

```
# A tibble: 638 x 5
    ..1 mode
                      year month avg_trips
  <dbl> <chr>
                         <chr> <chr> <chr>
      1 All Modes by Qtr 2007
                              01
                                    NA
      2 Boat
                         2007 01
      3 Bus
                         2007 01
                                    335.819
                      2007
      4 Commuter Rail
                                    142.2
                     2007 01
                                    435.294
      5 Heavy Rail
                     2007 01
      6 Light Rail
                                    227.231
      7 Pct Chg / Yr
                         2007 01
                                    0.02
      8 Private Bus
                                    4.772
                         2007 01
      9 RIDE
                         2007 01
                                    4.9
     10 Trackless Trolley 2007 01
                                    12.757
# ... with 628 more rows
```

- year, month, avg_trips are all character variables
- There is an odd column named . . 1
- The rows with "All Modes by Qtr" and "TOTAL" aren't necessary

```
mbta2 %>%
  mutate(
    year = parse_number(year),
    month = parse_number(month),
    avg_trips = parse_number(avg_trips)
)
```

```
# A tibble: 638 x 5
    ..1 mode
                           year month avg_trips
  <dbl> <chr>
                          <dbl> <dbl>
                                          <dbl>
      1 All Modes by Qtr
                           2007
                                          NA
                           2007
                                          4
      2 Boat
      3 Bus
                           2007
                                         336.
      4 Commuter Rail
                           2007
                                         142.
      5 Heavy Rail
                           2007
                                         435.
      6 Light Rail
                                         227.
                           2007
      7 Pct Chg / Yr
                           2007
                                           0.02
      8 Private Bus
                           2007
                                         4.77
      9 RIDE
                           2007
                                         4.9
     10 Trackless Trolley
                           2007
                                          12.8
# ... with 628 more rows
```

```
mbta2 %>%
  mutate(
    year = parse_number(year),
    month = parse_number(month),
    avg_trips = parse_number(avg_trips)
    ) %>%
    select(-1)
```

```
# A tibble: 638 x 4
  mode
                     year month avg_trips
  <chr>
                    <dbl> <dbl>
                                    <db1>
 1 All Modes by Qtr
                     2007
                                    NA
 2 Boat
                     2007
                                   4
 3 Bus
                     2007
                                   336.
 4 Commuter Rail
                     2007
                                   142.
 5 Heavy Rail
                                   435.
                     2007
6 Light Rail
                     2007
                                   227.
 7 Pct Chg / Yr
                     2007
                                   0.02
 8 Private Bus
                     2007
                                  4.77
9 RIDE
                     2007
                                  4.9
10 Trackless Trolley 2007
                                   12.8
# ... with 628 more rows
```

```
mbta2 %>%
  mutate(
    year = parse_number(year),
    month = parse_number(month),
    avg_trips = parse_number(avg_trips)
) %>%
  select(-1) %>%
  filter(mode != 'TOTAL', mode != "All Modes by Qtr")
```

```
# A tibble: 522 x 4
  mode
                     year month avg_trips
  <chr>
                    <dbl> <dbl>
                                    <dbl>
 1 Boat
                     2007
                                     4
 2 Bus
                                   336.
                     2007
 3 Commuter Rail
                     2007
                                   142.
 4 Heavy Rail
                                   435.
                     2007
 5 Light Rail
                                   227.
                     2007
 6 Pct Chg / Yr
                     2007
                                   0.02
 7 Private Bus
                     2007
                                   4.77
 8 RIDE
                     2007
                                   4.9
 9 Trackless Trolley 2007
                                   12.8
10 Boat
                     2007
                                     3.6
# ... with 512 more rows
```

Other cleaning tasks

- 1. distinct() keeps the unique (non-duplicate) rows of a dataset. Usage: dataset %>% distinct()
- 2. If you want to keep only rows with complete data, you can invoke drop_na. Usage: dataset %>% drop_na(). You can modify drop_na by specifying variables from which you want to drop the missing values.
- 3. If you want to convert a value to missing (commonly 99 is used for missing data), then you can use replace_na within mutate to change to missing values on a column-by-column basis. Usage: dataset %>% mutate(var1 = na_if(var1, 99))

Cleaning Excel data

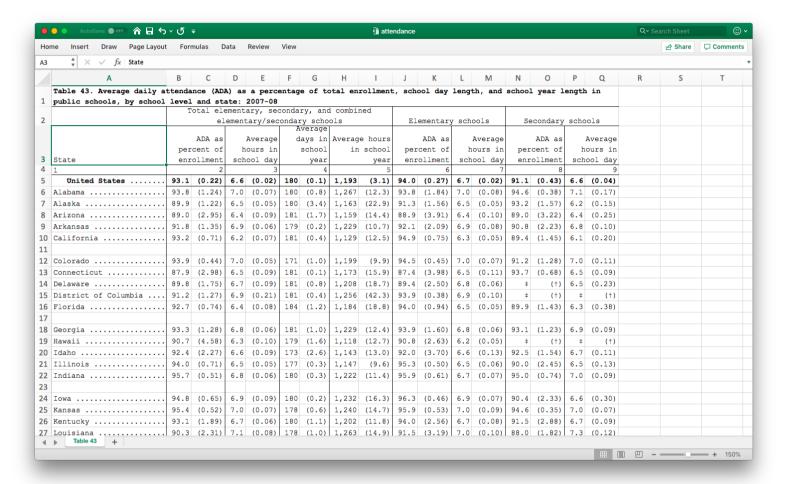
Excel is used as a visual medium

- Tables created to look good rather than being tidy or computer-friendly
 - Color being used to denote values of some variables
 - Multiple lines of headers
 - Multiple rows with variables
 - Typos leading to numeric variables become character

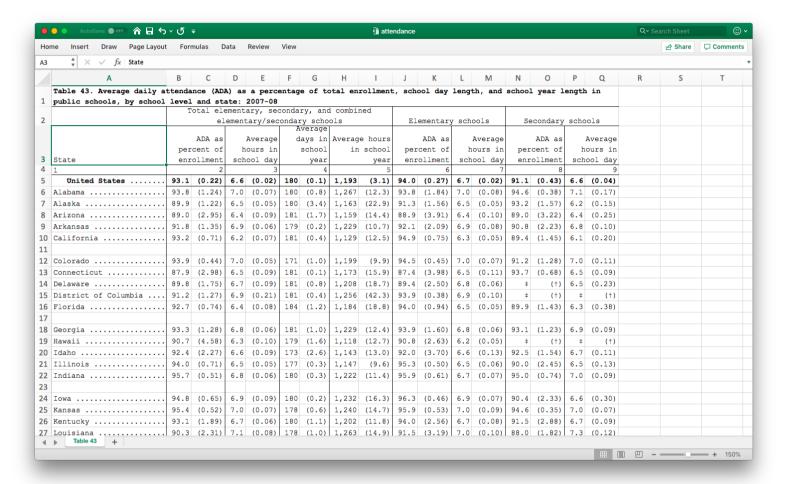
Excel is not reproducible, prone to mistakes by click

- Two special cases of Excel errors in the press
 - Duke cancer scandal with Dr. Anil Potti's group
 - Reinhart & Rogoff models for economic growth
- 35% of datasets in Nature (the journal) have Excel errors (The Economist, 2016)
 - A gene named 1MAR is entered in Excel. What does it become?

- Real data lies in the paired statistics and standard error columns
- The headers are basically different groupings and categories and should be variables



- import fails horribly
- Two packages, tidyxl and unpivotr, by Duncan Garmonsway, save the day



```
library(tidyxl)
dataset1 <- xlsx_cells('Data/FSI/attendance.xlsx')
dataset1</pre>
```

```
# A tibble: 1.173 x 21
                         col is_blank data_type error logical numeric
   sheet address
                   row
                 <int> <int> <lgl>
                                       <chr>
                                                 <chr> <lgl>
                                                                  <dbl>
   <chr> <chr>
                                       character <NA> NA
 1 Tabl... A1
                            1 FALSE
                                                                     NA
 2 Tabl... B1
                            2 TRUE
                                                 <NA> NA
                                       blank
                                                                     NA
 3 Tabl... C1
                            3 TRUE
                                       blank
                                                 <NA> NA
                                                                     NA
 4 Tabl... D1
                            4 TRUE
                                       blank
                                                 <NA> NA
                                                                     NA
 5 Tabl... E1
                            5 TRUE
                                       blank
                                                 <NA> NA
                                                                     NA
 6 Tabl... F1
                            6 TRUE
                                       blank
                                                 <NA> NA
                                                                     NA
 7 Tabl... G1
                           7 TRUE
                                       blank
                                                 <NA> NA
                                                                     NA
 8 Tabl... H1
                            8 TRUE
                                       blank
                                                 <NA> NA
                                                                     NA
 9 Tabl... I1
                            9 TRUE
                                       blank
                                                 <NA> NA
                                                                     NA
10 Tabl... J1
                           10 TRUE
                                       blank
                                                 <NA> NA
                                                                     NA
# ... with 1,163 more rows, and 12 more variables: date <dttm>,
    character <chr>, character_formatted <list>, formula <chr>,
  is_array <lgl>, formula_ref <chr>, formula_group <int>, comment <chr>,
   height <dbl>, width <dbl>, style_format <chr>, local_format_id <int>
```

- This grabs a bunch of meta-data about the Excel entries, including color and formatting features
- The data has been blown up on a cell-by-cell basis
- Use tidyverse tools to fix this? Nope. unpivotr is more powerful in this case.

```
library(unpivotr)
Warning: package 'unpivotr' was built under R version 3.5.2
dataset1 %>%
  filter(row != 1, row != 4, row < 65) %>%
  behead('N', tophead) %>%
  behead('N', head2) %>%
  behead('W', State) %>%
  select(row, col, data_type, numeric, tophead, head2, State)
# A tibble: 960 x 7
                          numeric tophead
     row
           col data_type
                                                           head2
                                                                     State
                             <dbl> <chr>
   <int> <int> <chr>
                                                           <chr>
                                                                     <chr>
             2 numeric
                           9.31e+1 Total elementary, se... ADA as p...
                                                                         Unit...
                                                                         Unit...
             3 numeric
                           2.19e-1 <NA>
                                                           <NA>
             4 numeric
                           6.64e+0 <NA>
                                                                         Unit...
                                                           Average ...
             5 numeric
                           1.76e-2 <NA>
                                                           <NA>
                                                                         Unit...
                         1.80e+2 <NA>
             6 numeric
                                                                         Unit...
                                                           Average ...
             7 numeric
                           1.43e-1 <NA>
                                                           <NA>
                                                                         Unit...
             8 numeric
                         1.19e+3 <NA>
                                                                         Unit...
                                                          Average ...
                           3.09e+0 <NA>
             9 numeric
                                                           <NA>
                                                                         Unit...
                         9.40e+1 Elementary schools
                                                          ADA as p...
          10 numeric
                                                                         Unit...
            11 numeric
                           2.69e-1 < NA >
                                                                         Unit...
                                                           <NA>
# ... with 950 more rows
```

• Pull off the two headers first with behead. Tell the function what direction (N, W, S, E or angles) the header is sitting in relation to the data

```
library(unpivotr)
dataset1 %>%
  filter(row != 1, row != 4, row < 65) %>%
  behead('N', tophead) %>%
  behead('N', head2) %>%
  behead('W', State) %>%
  select(row, col, data_type, numeric, tophead, head2, State) %>%
  mutate(header = ifelse(col %% 2 == 0, 'stats','se'))
```

```
# A tibble: 960 x 8
          col data_type numeric tophead
                                                 head2
                                                                  header
                                                          State
  <int> <int> <chr>
                       <dbl> <chr>
                                                 <chr>
                                                          <chr>
                                                                  <chr>
                                                             Uni... stats
            2 numeric
                      9.31e+1 Total elementar... ADA as ...
            3 numeric
                      2.19e-1 <NA>
                                                 <NA>
                                                             Uni… se
         4 numeric
                      6.64e+0 <NA>
                                                 Average... "
                                                             Uni... stats
            5 numeric
                       1.76e-2 <NA>
                                                 <NA>
                                                             Uni… se
           6 numeric
                                                 Average... "
                       1.80e+2 <NA>
                                                             Uni... stats
         7 numeric
                       1.43e-1 <NA>
                                                 <NA>
                                                             Uni… se
      5 8 numeric
                       1.19e+3 <NA>
                                                 Average... "
                                                             Uni... stats
        9 numeric
                      3.09e+0 <NA>
                                                 <NA>
                                                             Uni… se
      5 10 numeric
                       9.40e+1 Elementary scho… ADA as … "
                                                             Uni... stats
           11 numeric
                         2.69e-1 < NA >
                                                             Uni… se
                                                 <NA>
# ... with 950 more rows
```

- even columns are stats, odd columns are standard errors
- %% gives the remainder when left side is divided by right side

```
library(unpivotr)
dataset1 %>%
  filter(row != 1, row != 4, row < 65) %>%
  behead('N', tophead) %>%
  behead('N', head2) %>%
  behead('W', State) %>%
  select(row, col, data_type, numeric, tophead, head2, State) %>%
  mutate(header = ifelse(col %% 2 == 0, 'stats', 'se')) %>%
  fill(tophead) %>% fill(head2)
```

```
# A tibble: 960 x 8
               col data_type numeric tophead
                                                                      head2
                                                                                         State
                                                                                                      header
    <int> <int> <chr>
                                         <dbl> <chr>
                                                                           <chr>
                                                                                         <chr> <chr>
                  2 numeric 9.31e+1 Total elementar... ADA as ... "
                                                                                               Uni… stats
              3 numeric 2.19e-1 Total elementar… ADA as … "
                                                                                               Uni… se
       4 numeric
5 1.76e-2 Total elementar... Average... "
5 numeric
6.64e+0 Total elementar... Average... "
5 numeric
1.76e-2 Total elementar... Average... "
5 numeric
7 numeric
1.43e-1 Total elementar... Average... "
5 numeric
8 numeric
1.19e+3 Total elementar... Average... "
5 numeric
9 numeric
1.09e+0 Total elementar... Average... "
                                                                                               Uni... stats
                                                                                               Uni… se
                                                                                               Uni... stats
                                                                                               Uni… se
                                                                                               Uni... stats
                                                                                               Uni… se
      5 10 numeric
                                   9.40e+1 Elementary scho… ADA as … "
                                                                                               Uni... stats
                                      2.69e-1 Elementary scho… ADA as …
                11 numeric
                                                                                               Uni… se
# ... with 950 more rows
```

- column headers spanned several columns visually, but rested in left-most column internally
- used last value carried forward to fill in the other columns

```
library(unpivotr)
tidy_dataset <- dataset1 %>%
 filter(row != 1, row != 4, row < 65) %>%
 behead('N', tophead) %>%
 behead('N', head2) %>%
 behead('W', State) %>%
  select(row, col, data_type, numeric, tophead, head2, State) %>%
 mutate(header = ifelse(col %% 2 == 0, 'stats', 'se')) %>%
 fill(tophead) %>% fill(head2) %>%
 select(-col) %>%
 spatter(header, numeric) %>%
 select(-row)
tidv_dataset
# A tibble: 480 x 5
                                  head2
                                                 State
  tophead
                                                                  se stats
   <chr>
                                                               <dbl> <dbl>
                                  <chr>
                                                 <chr>
                                  ADA as percen... " United ... 0.269 9.40e1
 1 Elementary schools
                                  Average hours... " United ... 0.0160 6.66e0
 2 Elementary schools
                                 ADA as percen... " United ... 0.432 9.11e1
 3 Secondary schools
                                  Average hours... "
 4 Secondary schools
                                                     United ... 0.0403 6.59e0
 5 Total elementary, secondary, ... ADA as percen... "
                                                     United ... 0.219 9.31e1
 6 Total elementary, secondary, ... Average days ... "
                                                     United ... 0.143 1.80e2
7 Total elementary, secondary, ... Average hours... "
                                                     United ... 0.0176 6.64e0
 8 Total elementary, secondary, ... Average hours... "
                                                     United ... 3.09 1.19e3
 9 Elementary schools
                        ADA as percen... Alabama ..... 1.84 9.38e1
10 Elementary schools
                                  Average hours... Alabama ..... 0.0759 7.04e0
# ... with 470 more rows
```

• spatter works like spread, but is more robust for this kind of weird data

```
tidy_dataset <- tidy_dataset %>%
  mutate(State = str_remove(State, '\\.+')) %>%
  mutate(State = str_trim(State))
tidy_dataset
```

```
# A tibble: 480 x 5
   tophead
                                      head2
                                                       State
                                                                     se stats
                                                                  <dbl> <dbl>
   <chr>
                                      <chr>
                                                       <chr>
                                      ADA as percent... United ... 0.269 9.40e1
 1 Elementary schools
 2 Elementary schools
                                      Average hours ... United ... 0.0160 6.66e0
 3 Secondary schools
                                      ADA as percent... United ... 0.432 9.11e1
 4 Secondary schools
                                      Average hours ... United ... 0.0403 6.59e0
 5 Total elementary, secondary, and... ADA as percent... United ... 0.219 9.31e1
6 Total elementary, secondary, and... Average days i... United ... 0.143 1.80e2
 7 Total elementary, secondary, and... Average hours ... United ... 0.0176 6.64e0
 8 Total elementary, secondary, and... Average hours ... United ... 3.09
 9 Elementary schools
                                      ADA as percent... Alabama 1.84 9.38e1
10 Elementary schools
                                      Average hours ... Alabama 0.0759 7.04e0
# ... with 470 more rows
```

- We're using a **regular expression** to identify and remove all the dots
 - Rich tool for text searching
- Next we trim away any white space that is around the string

Save the data

```
saveRDS(tidy_dataset, file = 'Data/FSI/schools.rds')
```

The RDS format is an open standard and a fast way to store and retrieve datasets in R

What about being colorful?

```
library(tidyxl)
library(unpivotr)

dataset2 <- xlsx_cells('Data/FSI/classlist.xlsx')
formats <- xlsx_formats('Data/FSI/classlist.xlsx')</pre>
```

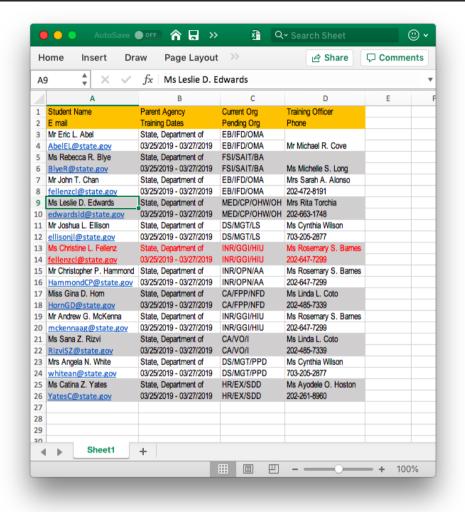
We need to grab the formats too, now

```
format_id <- dataset2$local_format_id
dataset2$font_color <- formats$local$font$color$rgb[f
dataset2$bg_color <- formats$local$fill$patternFill$f
unique(dataset2$font_color)

[1] "FF000000" "FF0563C1" "FFFF0000"

unique(dataset2$bg_color)

[1] "FFFFC000" NA "FFE7E6E6"</pre>
```



Grab the red rows

```
red_rows <- dataset2 %>% filter(font_color=='FFFF0000') %>%
    select(row, col, data_type, character) %>%
    mutate(row=2, col = 1:8)
headers <- dataset2 %>% filter(bg_color == 'FFFFC000') %>%
    select(row, col, data_type, character) %>%
    mutate(row = 1, col = 1:8)

bind_rows(headers, red_rows) %>%
    behead('N', header) %>%
    select(-col) %>%
    spatter(header) %>%
    select(-row)
```

There are really two datasets interwoven here

- The odd rows form one dataset
- The even rows form another dataset

We need to put these two datasets side-by-side

```
dat1 <- dataset2 %>%
  filter( row %% 2 == 1) %>% # odd rows
  behead('N', header) %>%
  mutate(row = (row+1)/2) # make the row numbers sequ

dat2 <- dataset2 %>%
  filter(row %% 2 == 0) %>% # even rows
  behead('N', header) %>%
  mutate(row = row/2) %>% # make row numbers sequenti
  mutate(col = col+4) # These will be the last 4 cols

tidy_dataset2 <-
  rbind(dat1, dat2) %>% # Put datsets on top of each
  select(row, data_type, numeric, character, header)
  spatter(header) %>%
  select(-row, -numeric)
```

```
dat1 <- dataset2 %>%
  filter( row %% 2 == 1) %>% # odd rows
  behead('N', header) %>%
  mutate(row = (row+1)/2) # make the row numbers sequ

dat2 <- dataset2 %>%
  filter(row %% 2 == 0) %>% # even rows
  behead('N', header) %>%
  mutate(row = row/2) %>% # make row numbers sequenti
  mutate(col = col+4) # These will be the last 4 cols

tidy_dataset2 <-
  rbind(dat1, dat2) %>% # Put datsets on top of each
  select(row, data_type, numeric, character, header)
  spatter(header) %>%
  select(-row, -numeric)
```

```
# A tibble: 12 x 8
   `Current Org` `E mail` `Parent Agency` `Pending Or
   <chr>
                  <chr>
                           <chr>
                                             <chr>
 1 EB/IFD/OMA
                  AbelEL@... State, Departm... EB/IFD/OMA
                  BlyeR@s... State, Departm... FSI/SAIT/BA
 2 FSI/SAIT/BA
                  fellenz... State, Departm... EB/IFD/OMA
 3 EB/IFD/OMA
 4 MED/CP/OHW/OH
                  edwards... State, Departm... MED/CP/OHW/
                  ellison... State, Departm... DS/MGT/LS
 5 DS/MGT/LS
 6 INR/GGI/HIU
                  fellenz... State, Departm... INR/GGI/HIU
 7 INR/OPN/AA
                  Hammond... State, Departm... INR/OPN/AA
                  HornGD@... State, Departm... CA/FPP/NFD
 8 CA/FPP/NFD
 9 INR/GGI/HIU
                  mckenna... State, Departm... INR/GGI/HIU
                  RizviSZ... State, Departm... CA/VO/I
10 CA/VO/I
11 DS/MGT/PPD
                  whitean... State, Departm... DS/MGT/PPD
12 HR/EX/SDD
                  YatesC@... State, Departm... HR/EX/SDD
# ... with 3 more variables: `Student Name` <chr>, `Tra
    `Training Officer` <chr>
```

```
tidy_dataset2 <- tidy_dataset2 %>%
   set_names(make.names(names(.))) %>%
   select(Student.Name, everything())
```

```
# A tibble: 12 x 8
   Student.Name Current.Org E.mail Parent.Agency Pending.Org Phone
   <chr>
                 <chr>
                               <chr> <chr>
                                                      <chr>
                                                                   <chr>
 1 Mr Eric L. ... EB/IFD/OMA
                              AbelE... State, Depar... EB/IFD/OMA
                                                                   <NA>
 2 Ms Rebecca ... FSI/SAIT/BA BlyeR... State, Depar... FSI/SAIT/BA <NA>
 3 Mr John T. ... EB/IFD/OMA
                             felle... State, Depar... EB/IFD/OMA
 4 Ms Leslie D... MED/CP/OHW... edwar... State, Depar... MED/CP/OHW... 202-...
 5 Mr Joshua L... DS/MGT/LS
                               ellis... State, Depar... DS/MGT/LS
 6 Ms Christin... INR/GGI/HIU felle... State, Depar... INR/GGI/HIU 202-...
 7 Mr Christop... INR/OPN/AA
                              Hammo... State, Depar... INR/OPN/AA
                                                                  202-...
 8 Miss Gina D., CA/FPP/NFD
                              HornG... State, Depar... CA/FPP/NFD
                                                                  202-...
 9 Mr Andrew G... INR/GGI/HIU mcken... State, Depar... INR/GGI/HIU 202-...
10 Ms Sana Z. ... CA/VO/I
                               Rizvi... State, Depar... CA/VO/I
                                                                   202-...
11 Mrs Angela ... DS/MGT/PPD
                              white... State, Depar... DS/MGT/PPD 703-...
12 Ms Catina Z... HR/EX/SDD
                               Yates... State, Depar... HR/EX/SDD
                                                                   202-...
# ... with 2 more variables: Training.Dates <chr>, Training.Officer <chr>
```