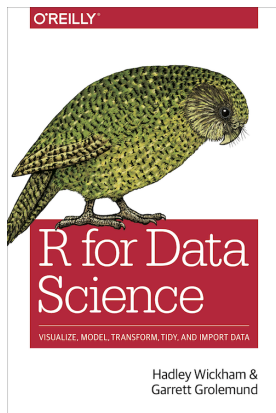# Socio-Informatics 348
## Data Tidying

Dr Lisa Martin

Department of Information Science
Stellenbosch University
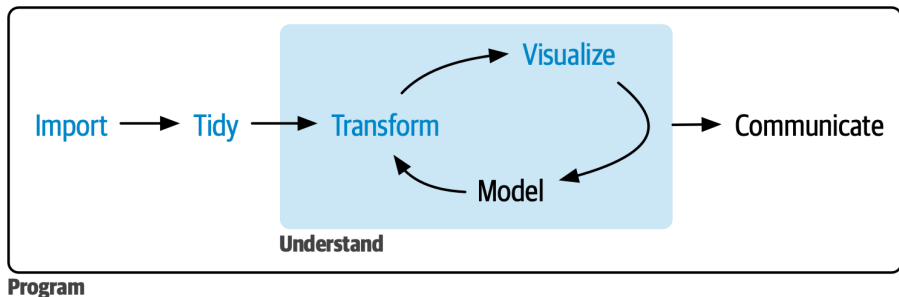
# Today's Reading



*R for Data Science, Wholegame, Data Tidying*

# Let's refer back to the 'wholegame'

From *R for Data Science*, Introduction:

# Data Structure

**Rules**

- Each variable is a column.
- Each observation is a row, represents a single unit at a point in time.
- Each cell is a single value.



variables      observations      values

**Advantages**

- R's vectorized nature
- Consistency makes it easier to learn and use

# Real-world data is messy

- Not often in tidy format
- Often compiled for a specific purpose, not for your research

# Lengthening Data

In some cases, there may be extra data captured in column headings

```
billboard
#> # A tibble: 317 × 79
#>   artist       track               date.entered   wk1   wk2   wk3   wk4   wk5
#>   <chr>        <chr>               <date>       <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1 2 Pac        Baby Don't Cry (Ke… 2000-02-26      87    82    72    77    87
#> 2 2Ge+her      The Hardest Part O… 2000-09-02      91    87    92    NA    NA
#> 3 3 Doors Down Kryptonite          2000-04-08      81    70    68    67    66
#> 4 3 Doors Down Loser               2000-10-21      76    76    72    69    67
#> 5 504 Boyz     Wobble Wobble       2000-04-15      57    34    25    17    17
#> 6 98^0         Give Me Just One N… 2000-08-19      51    39    34    26    26
#> # i 311 more rows
#> # i 71 more variables: wk6 <dbl>, wk7 <dbl>, wk8 <dbl>, wk9 <dbl>, …
```

# Lengthening Data

Here, we have time information in the column headings, which is not ideal

- We would prefer to have a single column for 'week'
- Think of how that might help with visualising or transorming the data

```
billboard
#> # A tibble: 317 × 79
#>    artist       track                date.entered   wk1   wk2   wk3   wk4   wk5
#>    <chr>        <chr>                <date>       <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1 2 Pac        Baby Don't Cry (Ke… 2000-02-26      87    82    72    77    87
#> 2 2Ge+her      The Hardest Part O… 2000-09-02      91    87    92    NA    NA
#> 3 3 Doors Down Kryptonite          2000-04-08      81    70    68    67    66
#> 4 3 Doors Down Loser               2000-10-21      76    76    72    69    67
#> 5 504 Boyz     Wobble Wobble       2000-04-15      57    34    25    17    17
#> 6 98^0         Give Me Just One N… 2000-08-19      51    39    34    26    26
#> # i 311 more rows
#> # i 71 more variables: wk6 <dbl>, wk7 <dbl>, wk8 <dbl>, wk9 <dbl>, …
```

# Lengthening Data

**pivot_longer**

```r
billboard |>
  pivot_longer(
    cols = starts_with("wk"),
    names_to = "week",
    values_to = "rank"
  )
#> # A tibble: 24,092 × 5
#>    artist track                    date.entered week   rank
#>    <chr>  <chr>                    <date>       <chr> <dbl>
#>  1 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk1      87
#>  2 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk2      82
#>  3 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk3      72
#>  4 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk4      77
#>  5 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk5      87
#>  6 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk6      94
#>  7 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk7      99
#>  8 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk8      NA
#>  9 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk9      NA
#> 10 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk10     NA
#> # i 24,082 more rows
```

# Lengthening Data

**pivot_longer**



*Adapted from Visuals by Garrick Aden-Buie*

- Think of how you would plot this new data

# Lengthening Data

**pivot_longer**

```r
billboard |>
  pivot_longer(
    cols = starts_with("wk"),
    names_to = "week",
    values_to = "rank",
    values_drop_na = TRUE
  )
#> # A tibble: 5,307 × 5
#>   artist track                    date.entered week   rank
#>   <chr>  <chr>                    <date>       <chr> <dbl>
#> 1 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk1      87
#> 2 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk2      82
#> 3 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk3      72
#> 4 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk4      77
#> 5 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk5      87
#> 6 2 Pac  Baby Don't Cry (Keep...  2000-02-26   wk6      94
#> # i 5,301 more rows
```

# Lengthening Data

**pivot_longer**

- Each column is a variable
- Each cell is a single value
- But now we have multiple rows for each 'unit' or 'object'
- Each row is still a single observation in time

**Multiple variables in a single column name?**

- Collected by the World Health Organisation
- Records information about tuberculosis diagnoses

```
who2
#> # A tibble: 7,240 × 58
#>    country     year sp_m_014 sp_m_1524 sp_m_2534 sp_m_3544 sp_m_4554
#>    <chr>      <dbl>    <dbl>     <dbl>     <dbl>     <dbl>     <dbl>
#> 1 Afghanistan 1980       NA        NA        NA        NA        NA
#> 2 Afghanistan 1981       NA        NA        NA        NA        NA
#> 3 Afghanistan 1982       NA        NA        NA        NA        NA
#> 4 Afghanistan 1983       NA        NA        NA        NA        NA
#> 5 Afghanistan 1984       NA        NA        NA        NA        NA
#> 6 Afghanistan 1985       NA        NA        NA        NA        NA
#> # i 7,234 more rows
#> # i 51 more variables: sp_m_5564 <dbl>, sp_m_65 <dbl>, sp_f_014 <dbl>, …
```

# Lengthening Data

**Multiple variables in a single column name?**

```r
who2 |>
  pivot_longer(
    cols = !(country:year),
    names_to = c("diagnosis", "gender", "age"),
    names_sep = "_",
    values_to = "count"
  )
#> # A tibble: 405,440 × 6
#>   country     year diagnosis gender age   count
#>   <chr>      <dbl> <chr>     <chr>  <chr> <dbl>
#> 1 Afghanistan 1980 sp        m      014      NA
#> 2 Afghanistan 1980 sp        m      1524     NA
#> 3 Afghanistan 1980 sp        m      2534     NA
#> 4 Afghanistan 1980 sp        m      3544     NA
#> 5 Afghanistan 1980 sp        m      4554     NA
#> 6 Afghanistan 1980 sp        m      5564     NA
#> # i 405,434 more rows
```

# Widening Data

- Needing to widen data is less common
- `cms_patient_experience` from the Centers of Medicare and Medicaid services
- Data about patient experiences

```
cms_patient_experience
#> # A tibble: 500 × 5
#>   org_pac_id org_nm                     measure_cd    measure_title  prf_rate
#>   <chr>      <chr>                      <chr>         <chr>             <dbl>
#> 1 0446157747 USC CARE MEDICAL GROUP INC CAHPS_GRP_1   CAHPS for MIPS…      63
#> 2 0446157747 USC CARE MEDICAL GROUP INC CAHPS_GRP_2   CAHPS for MIPS…      87
#> 3 0446157747 USC CARE MEDICAL GROUP INC CAHPS_GRP_3   CAHPS for MIPS…      86
#> 4 0446157747 USC CARE MEDICAL GROUP INC CAHPS_GRP_5   CAHPS for MIPS…      57
#> 5 0446157747 USC CARE MEDICAL GROUP INC CAHPS_GRP_8   CAHPS for MIPS…      85
#> 6 0446157747 USC CARE MEDICAL GROUP INC CAHPS_GRP_12  CAHPS for MIPS…      24
#> # i 494 more rows
```

# Widening Data



Adapted from Visuals by Garrick Aden-Buie

# Widening Data

**pivot_wider**

```
cms_patient_experience |>
  pivot_wider(
    id_cols = starts_with("org"),  ⟵ unique identifier / unit for each row
    names_from = measure_cd,
    values_from = prf_rate
  )
#> # A tibble: 95 × 8
#>   org_pac_id org_nm      CAHPS_GRP_1 CAHPS_GRP_2 CAHPS_GRP_3 CAHPS_GRP_5
#>   <chr>      <chr>             <dbl>       <dbl>       <dbl>       <dbl>
#> 1 0446157747 USC CARE MEDICA…     63          87          86          57
#> 2 0446162697 ASSOCIATION OF …     59          85          83          63
#> 3 0547164295 BEAVER MEDICAL …     49          NA          75          44
#> 4 0749333730 CAPE PHYSICIANS…     67          84          85          65
#> 5 0840104360 ALLIANCE PHYSIC…     66          87          87          64
#> 6 0840109864 REX HOSPITAL INC     73          87          84          67
#> # i 89 more rows
#> # i 2 more variables: CAHPS_GRP_8 <dbl>, CAHPS_GRP_12 <dbl>
```