

12.6.1

Han Nguyen

February 25, 2018

```
library(tidyverse)
```

```
## Warning: package 'tidyverse' was built under R version 3.4.3
## -- Attaching packages ----- tidyverse 1.2.1 --
## v ggplot2 2.2.1    v purrr  0.2.4
## v tibble  1.4.2    v dplyr  0.7.4
## v tidyr   0.8.0    v stringr 1.2.0
## v readr   1.1.1    v forcats 0.2.0

## Warning: package 'ggplot2' was built under R version 3.4.3
## Warning: package 'tibble' was built under R version 3.4.3
## Warning: package 'tidyr' was built under R version 3.4.3
## Warning: package 'readr' was built under R version 3.4.3
## Warning: package 'purrr' was built under R version 3.4.3
## Warning: package 'dplyr' was built under R version 3.4.3
## Warning: package 'stringr' was built under R version 3.4.2
## Warning: package 'forcats' was built under R version 3.4.3

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
library(tidyr)
```

```
who
```

```
## # A tibble: 7,240 x 60
##   country    iso2 iso3  year new_sp_m014 new_sp_m1524 new_sp_m2534
##   <chr>      <chr> <chr> <int>      <int>      <int>      <int>
## 1 Afghanistan AF    AFG   1980         NA         NA         NA
## 2 Afghanistan AF    AFG   1981         NA         NA         NA
## 3 Afghanistan AF    AFG   1982         NA         NA         NA
## 4 Afghanistan AF    AFG   1983         NA         NA         NA
## 5 Afghanistan AF    AFG   1984         NA         NA         NA
## 6 Afghanistan AF    AFG   1985         NA         NA         NA
## 7 Afghanistan AF    AFG   1986         NA         NA         NA
## 8 Afghanistan AF    AFG   1987         NA         NA         NA
## 9 Afghanistan AF    AFG   1988         NA         NA         NA
## 10 Afghanistan AF    AFG   1989         NA         NA         NA
## # ... with 7,230 more rows, and 53 more variables: new_sp_m3544 <int>,
## #   new_sp_m4554 <int>, new_sp_m5564 <int>, new_sp_m65 <int>,
## #   new_sp_f014 <int>, new_sp_f1524 <int>, new_sp_f2534 <int>,
## #   new_sp_f3544 <int>, new_sp_f4554 <int>, new_sp_f5564 <int>,
## #   new_sp_f65 <int>, new_sn_m014 <int>, new_sn_m1524 <int>,
## #   new_sn_m2534 <int>, new_sn_m3544 <int>, new_sn_m4554 <int>,
## #   new_sn_m5564 <int>, new_sn_m65 <int>, new_sn_f014 <int>,
```

```
## # new_sn_f1524 <int>, new_sn_f2534 <int>, new_sn_f3544 <int>,
## # new_sn_f4554 <int>, new_sn_f5564 <int>, new_sn_f65 <int>,
## # new_ep_m014 <int>, new_ep_m1524 <int>, new_ep_m2534 <int>,
## # new_ep_m3544 <int>, new_ep_m4554 <int>, new_ep_m5564 <int>,
## # new_ep_m65 <int>, new_ep_f014 <int>, new_ep_f1524 <int>,
## # new_ep_f2534 <int>, new_ep_f3544 <int>, new_ep_f4554 <int>,
## # new_ep_f5564 <int>, new_ep_f65 <int>, newrel_m014 <int>,
## # newrel_m1524 <int>, newrel_m2534 <int>, newrel_m3544 <int>,
## # newrel_m4554 <int>, newrel_m5564 <int>, newrel_m65 <int>,
## # newrel_f014 <int>, newrel_f1524 <int>, newrel_f2534 <int>,
## # newrel_f3544 <int>, newrel_f4554 <int>, newrel_f5564 <int>,
## # newrel_f65 <int>
```

```
who1 <- who %>%
  gather(new_sp_m014:newrel_f65, key = "key", value = "cases", na.rm = TRUE)
who1
```

```
## # A tibble: 76,046 x 6
##   country    iso2 iso3   year key      cases
##   * <chr>      <chr> <chr> <int> <chr>    <int>
## 1 Afghanistan AF    AFG   1997 new_sp_m014    0
## 2 Afghanistan AF    AFG   1998 new_sp_m014   30
## 3 Afghanistan AF    AFG   1999 new_sp_m014    8
## 4 Afghanistan AF    AFG   2000 new_sp_m014   52
## 5 Afghanistan AF    AFG   2001 new_sp_m014  129
## 6 Afghanistan AF    AFG   2002 new_sp_m014   90
## 7 Afghanistan AF    AFG   2003 new_sp_m014  127
## 8 Afghanistan AF    AFG   2004 new_sp_m014  139
## 9 Afghanistan AF    AFG   2005 new_sp_m014  151
## 10 Afghanistan AF    AFG   2006 new_sp_m014  193
## # ... with 76,036 more rows
```

```
who1 %>%
  count(key)
```

```
## # A tibble: 56 x 2
##   key      n
##   <chr>    <int>
## 1 new_ep_f014  1032
## 2 new_ep_f1524 1021
## 3 new_ep_f2534 1021
## 4 new_ep_f3544 1021
## 5 new_ep_f4554 1017
## 6 new_ep_f5564 1017
## 7 new_ep_f65   1014
## 8 new_ep_m014  1038
## 9 new_ep_m1524 1026
## 10 new_ep_m2534 1020
## # ... with 46 more rows
```

```
who2 <- who1 %>%
  mutate(key = stringr::str_replace(key, "newrel", "new_rel"))
```

```
## Warning: package 'bindrcpp' was built under R version 3.4.3
who2
```

```
## # A tibble: 76,046 x 6
##   country    iso2 iso3   year key      cases
##   <chr>      <chr> <chr> <int> <chr>    <int>
## 1 Afghanistan AF    AFG   1997 new_sp_m014    0
## 2 Afghanistan AF    AFG   1998 new_sp_m014   30
## 3 Afghanistan AF    AFG   1999 new_sp_m014    8
## 4 Afghanistan AF    AFG   2000 new_sp_m014   52
## 5 Afghanistan AF    AFG   2001 new_sp_m014  129
## 6 Afghanistan AF    AFG   2002 new_sp_m014   90
## 7 Afghanistan AF    AFG   2003 new_sp_m014  127
## 8 Afghanistan AF    AFG   2004 new_sp_m014  139
## 9 Afghanistan AF    AFG   2005 new_sp_m014  151
## 10 Afghanistan AF    AFG   2006 new_sp_m014  193
## # ... with 76,036 more rows
```

```
who3 <- who2 %>%
  separate(key, c("new", "type", "sexage"), sep = "_")
who3
```

```
## # A tibble: 76,046 x 8
##   country    iso2 iso3   year new   type sexage cases
##   <chr>      <chr> <chr> <int> <chr> <chr> <chr>    <int>
## 1 Afghanistan AF    AFG   1997 new   sp   m014    0
## 2 Afghanistan AF    AFG   1998 new   sp   m014   30
## 3 Afghanistan AF    AFG   1999 new   sp   m014    8
## 4 Afghanistan AF    AFG   2000 new   sp   m014   52
## 5 Afghanistan AF    AFG   2001 new   sp   m014  129
## 6 Afghanistan AF    AFG   2002 new   sp   m014   90
## 7 Afghanistan AF    AFG   2003 new   sp   m014  127
## 8 Afghanistan AF    AFG   2004 new   sp   m014  139
## 9 Afghanistan AF    AFG   2005 new   sp   m014  151
## 10 Afghanistan AF    AFG   2006 new   sp   m014  193
## # ... with 76,036 more rows
```

```
who3 %>%
  count(new)
```

```
## # A tibble: 1 x 2
##   new      n
##   <chr> <int>
## 1 new   76046
```

```
who4 <- who3 %>%
  select(-new, -iso2, -iso3)
```

```
who5 <- who4 %>%
  separate(sexage, c("sex", "age"), sep = 1)
who5
```

```
## # A tibble: 76,046 x 6
##   country    year type sex   age cases
##   <chr>      <int> <chr> <chr> <chr>    <int>
## 1 Afghanistan 1997 sp    m    014    0
## 2 Afghanistan 1998 sp    m    014   30
## 3 Afghanistan 1999 sp    m    014    8
## 4 Afghanistan 2000 sp    m    014   52
## 5 Afghanistan 2001 sp    m    014  129
```

```
## 6 Afghanistan 2002 sp m 014 90
## 7 Afghanistan 2003 sp m 014 127
## 8 Afghanistan 2004 sp m 014 139
## 9 Afghanistan 2005 sp m 014 151
## 10 Afghanistan 2006 sp m 014 193
## # ... with 76,036 more rows
```

Question 3

I claimed that iso2 and iso3 were redundant with country. Confirm this claim.

```
select(who3, country, iso2, iso3) %>%
  distinct() %>%
  group_by(country) %>%
  filter(n() > 1)
```

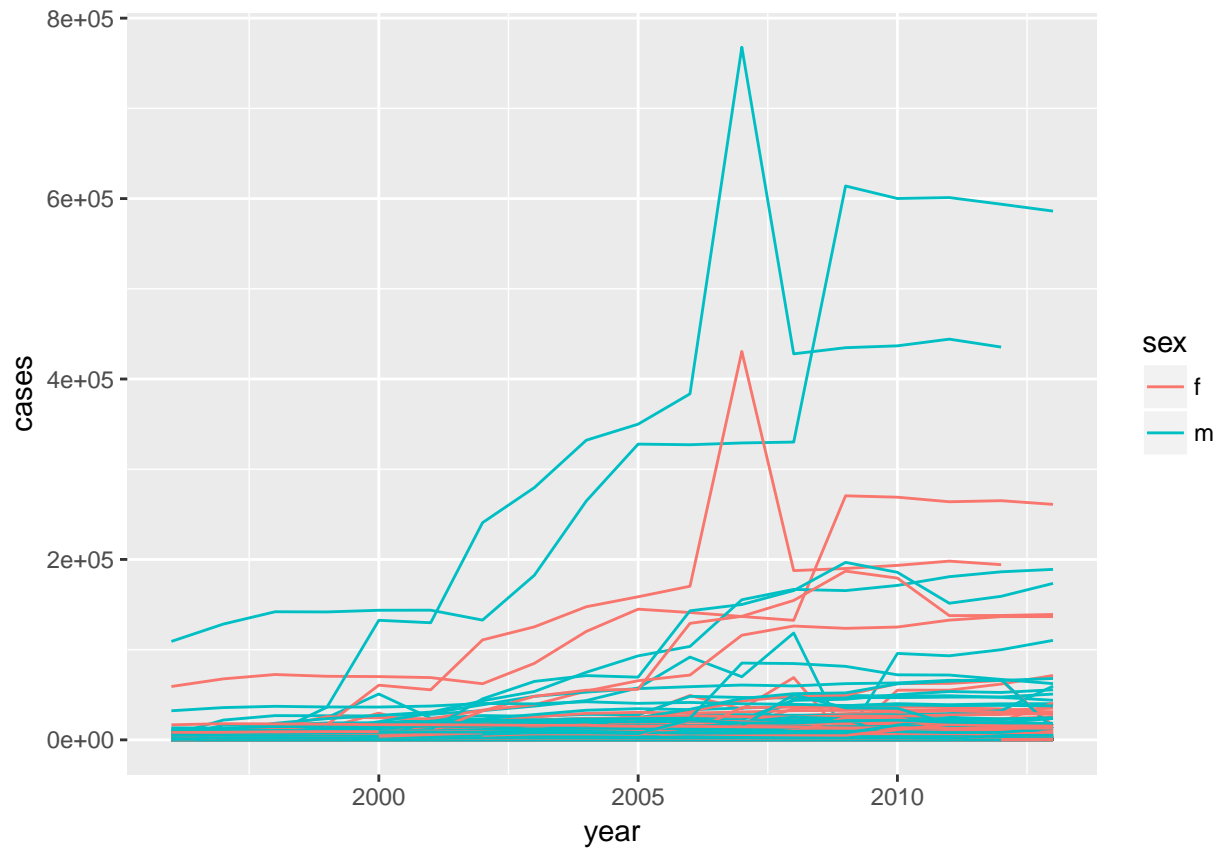
```
## # A tibble: 0 x 3
## # Groups:   country [0]
## # ... with 3 variables: country <chr>, iso2 <chr>, iso3 <chr>
```

Since this returns 0 distinct tibbles and groups, it shows that iso2 and iso3 are redundant with country.

Question 4

For each country, year, and sex compute the total number of cases of TB. Make an informative visualisation of the data.

```
who5 %>%
  group_by(country, year, sex) %>%
  filter(year > 1995) %>%
  summarise(cases = sum(cases)) %>%
  unite(country_sex, country, sex, remove = FALSE) %>%
  ggplot(aes(x = year, y = cases, group = country_sex, colour = sex)) +
  geom_line()
```



Countries were combined with sex because there were so many countries that a ggplot would not represent each one well.