

Assignment_4_part2_chap12

Anthea Yichen Li

February 20, 2018

1. In this case study I set `na.rm = TRUE` just to make it easier to check that we had the correct values. Is this reasonable? Think about how missing values are represented in this dataset. Are there implicit missing values? What's the difference between an NA and zero?

Perhaps? I would need to know more about the data generation process. There are zero's in the data, which means they may explicitly be indicating no cases.

```
who1 %>%  
  filter(cases == 0) %>%  
  nrow()
```

```
## [1] 11080
```

So it appears that either a country has all its values in a year as non-missing if the WHO collected data for that country, or all its values are non-missing. So it is okay to treat explicitly and implicitly missing values the same, and we don't lose any information by dropping them.

2. What happens if you neglect the `mutate()` step? (`mutate(key = stringr::str_replace(key, "newrel", "new_rel")`)

`separate` emits the warning "too few values", and if we check the rows for keys beginning with "newrel_", we see that `sexage` is messing, and `type = m014`.

```
who3a <- who1 %>%  
  separate(key, c("new", "type", "sexage"), sep = "_")
```

```
## Warning: Expected 3 pieces. Missing pieces filled with `NA` in 2580 rows  
## [73467, 73468, 73469, 73470, 73471, 73472, 73473, 73474, 73475, 73476,  
## 73477, 73478, 73479, 73480, 73481, 73482, 73483, 73484, 73485, 73486, ...].
```

```
filter(who3a, new == "newrel") %>% head()
```

```
## # A tibble: 6 x 8  
##   country      iso2 iso3   year new    type  sexage cases  
##   <chr>      <chr> <chr> <int> <chr>  <chr> <chr>  <int>  
## 1 Afghanistan AF    AFG   2013 newrel m014  <NA>    1705  
## 2 Albania    AL    ALB   2013 newrel m014  <NA>     14  
## 3 Algeria    DZ    DZA   2013 newrel m014  <NA>     25  
## 4 Andorra    AD    AND   2013 newrel m014  <NA>      0  
## 5 Angola     AO    AGO   2013 newrel m014  <NA>    486  
## 6 Anguilla   AI    AIA   2013 newrel m014  <NA>      0
```

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>
```

4. For each country, year, and sex compute the total number of cases of TB. Make an informative visualization 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()
```

