

Week 4 Exercises

Steven Simonsen

April 2, 2024

Please complete all exercises below. You may use any library that we have covered in class. The data we will be using comes from the tidyr package, so you must use that.

- 1) Examine the who and population data sets that come with the tidyr library. the who data is not tidy, you will need to reshape the new_sp_m014 to newrel_f65 columns to long format retaining country, iso2, iso3, and year. The data in the columns you are reshaping contains patterns described in the details section below. You will need to assign three columns: diagnosis, gender, and age to the patterns described in the details.

Your tidy data should look like the following: country iso2 iso3 year diagnosis gender age count
1 Afghanistan AF AFG 1980 sp m 014 NA 2 Afghanistan AF AFG 1980 sp m 1524 NA 3 Afghanistan AF AFG 1980 sp m 2534 NA 4 Afghanistan AF AFG 1980 sp m 3544 NA 5 Afghanistan AF AFG 1980 sp m 4554 NA 6 Afghanistan AF AFG 1980 sp m 5564 NA

Details The data uses the original codes given by the World Health Organization. The column names for columns five through 60 are made by combining new_ to a code for method of diagnosis (rel = relapse, sn = negative pulmonary smear, sp = positive pulmonary smear, ep = extrapulmonary) to a code for gender (f = female, m = male) to a code for age group (014 = 0-14 yrs of age, 1524 = 15-24 years of age, 2534 = 25 to 34 years of age, 3544 = 35 to 44 years of age, 4554 = 45 to 54 years of age, 5564 = 55 to 64 years of age, 65 = 65 years of age or older).

Note: use data(who) and data(population) to load the data into your environment. Use the arguments cols, names_to, names_pattern, and values_to. Your regex should be = ("new_?(.)_(.)(.)")

<https://tidyr.tidyverse.org/reference/who.html>

```
library(dplyr)
```

```
##  
## Attaching package: 'dplyr'  
  
## The following objects are masked from 'package:stats':  
##  
##   filter, lag  
  
## The following objects are masked from 'package:base':  
##  
##   intersect, setdiff, setequal, union
```

```
library(tidyr)  
library(ggplot2)
```

```
data(who)  
data("population")
```

```
#Pivot Long to match format above  
who_long_data <- who %>%
```

```

pivot_longer(cols = !country:year,
              names_to = c("diagnosis", "gender", "age"),
              names_pattern = ("new_?(.*)_(.)(.*)"),
              values_to = "count")

head(who_long_data)

```

```

## # A tibble: 6 x 8
##   country    iso2 iso3  year diagnosis gender age  count
##   <chr>      <chr> <chr> <dbl> <chr>      <chr> <chr> <dbl>
## 1 Afghanistan AF    AFG  1980 sp        m    014    NA
## 2 Afghanistan AF    AFG  1980 sp        m   1524    NA
## 3 Afghanistan AF    AFG  1980 sp        m   2534    NA
## 4 Afghanistan AF    AFG  1980 sp        m   3544    NA
## 5 Afghanistan AF    AFG  1980 sp        m   4554    NA
## 6 Afghanistan AF    AFG  1980 sp        m   5564    NA

```

- 2) There are two common keys between the data sets, with who as the left table, join the population data by country and year so that the population is available within the who dataset.

```

#Left join as to not lose any data in who dataset
left_join_who <- who_long_data %>%
  left_join(population, by=c('country', 'year'))

head(left_join_who)

```

```

## # A tibble: 6 x 9
##   country    iso2 iso3  year diagnosis gender age  count population
##   <chr>      <chr> <chr> <dbl> <chr>      <chr> <chr> <dbl>      <dbl>
## 1 Afghanistan AF    AFG  1980 sp        m    014    NA      NA
## 2 Afghanistan AF    AFG  1980 sp        m   1524    NA      NA
## 3 Afghanistan AF    AFG  1980 sp        m   2534    NA      NA
## 4 Afghanistan AF    AFG  1980 sp        m   3544    NA      NA
## 5 Afghanistan AF    AFG  1980 sp        m   4554    NA      NA
## 6 Afghanistan AF    AFG  1980 sp        m   5564    NA      NA

```

- 3) Split the age column into two columns, min age and max age. Notice that there is no character separator. Check the documentation with `?separate` to understand other ways to separate the age column. Keep in mind that 0 to 14 is coded as 014 (3 characters) and the other age groups are coded with 4 characters. 65 only has two characters, but we will ignore that until the next problem.

```

# Use sep=-2 to account for uneven character split
age_split_who <- left_join_who %>%
  separate(age, c("min_age", "max_age"),
           ,sep = -2
           )

head(age_split_who)

```

```

## # A tibble: 6 x 10
##   country    iso2 iso3  year diagnosis gender min_age max_age count population
##   <chr>      <chr> <chr> <dbl> <chr>      <chr> <chr> <chr> <dbl>      <dbl>
## 1 Afghanist~ AF    AFG  1980 sp        m     0     14      NA      NA
## 2 Afghanist~ AF    AFG  1980 sp        m    15     24      NA      NA
## 3 Afghanist~ AF    AFG  1980 sp        m    25     34      NA      NA
## 4 Afghanist~ AF    AFG  1980 sp        m    35     44      NA      NA
## 5 Afghanist~ AF    AFG  1980 sp        m    45     54      NA      NA

```

```
## 6 Afghanistan~ AF      AFG      1980 sp      m      55      64      NA      NA
```

- 4) Since we ignored the 65+ group in the previous problem we will fix it here. If you examine the data you will notice that 65 was placed into the max_age column and there is no value for min_age for those records. To fix this use mutate() in order to replace the blank value in the min_age column with the value from the max_age column and another mutate to replace the 65 in the max column with an Inf. Be sure to keep the variables as character vectors.

```
#kept min_age and max age as characters vectors as instructed with mutate.
age_split_who <- age_split_who %>%
  mutate(min_age = ifelse(min_age=="", 0, min_age)) %>%
  mutate(max_age = ifelse(max_age== 65, "Inf", max_age))

head(age_split_who)
```

```
## # A tibble: 6 x 10
##   country    iso2 iso3   year diagnosis gender min_age max_age count population
##   <chr>      <chr> <chr> <dbl> <chr>    <chr> <chr>  <chr>  <dbl>      <dbl>
## 1 Afghanistan~ AF      AFG    1980 sp      m      0      14      NA      NA
## 2 Afghanistan~ AF      AFG    1980 sp      m     15     24      NA      NA
## 3 Afghanistan~ AF      AFG    1980 sp      m     25     34      NA      NA
## 4 Afghanistan~ AF      AFG    1980 sp      m     35     44      NA      NA
## 5 Afghanistan~ AF      AFG    1980 sp      m     45     54      NA      NA
## 6 Afghanistan~ AF      AFG    1980 sp      m     55     64      NA      NA
```

- 5) Find the count per diagnosis for males and females.

See ?sum for a hint on resolving NA values.

```
#Use na.rm=TRUE to resolve NA values
gender_count <- age_split_who %>%
  group_by(gender, diagnosis) %>%
  summarize(count_per_diagnosis=sum(count, na.rm = TRUE))
```

```
## `summarise()` has grouped output by 'gender'. You can override using the
## `.groups` argument.
```

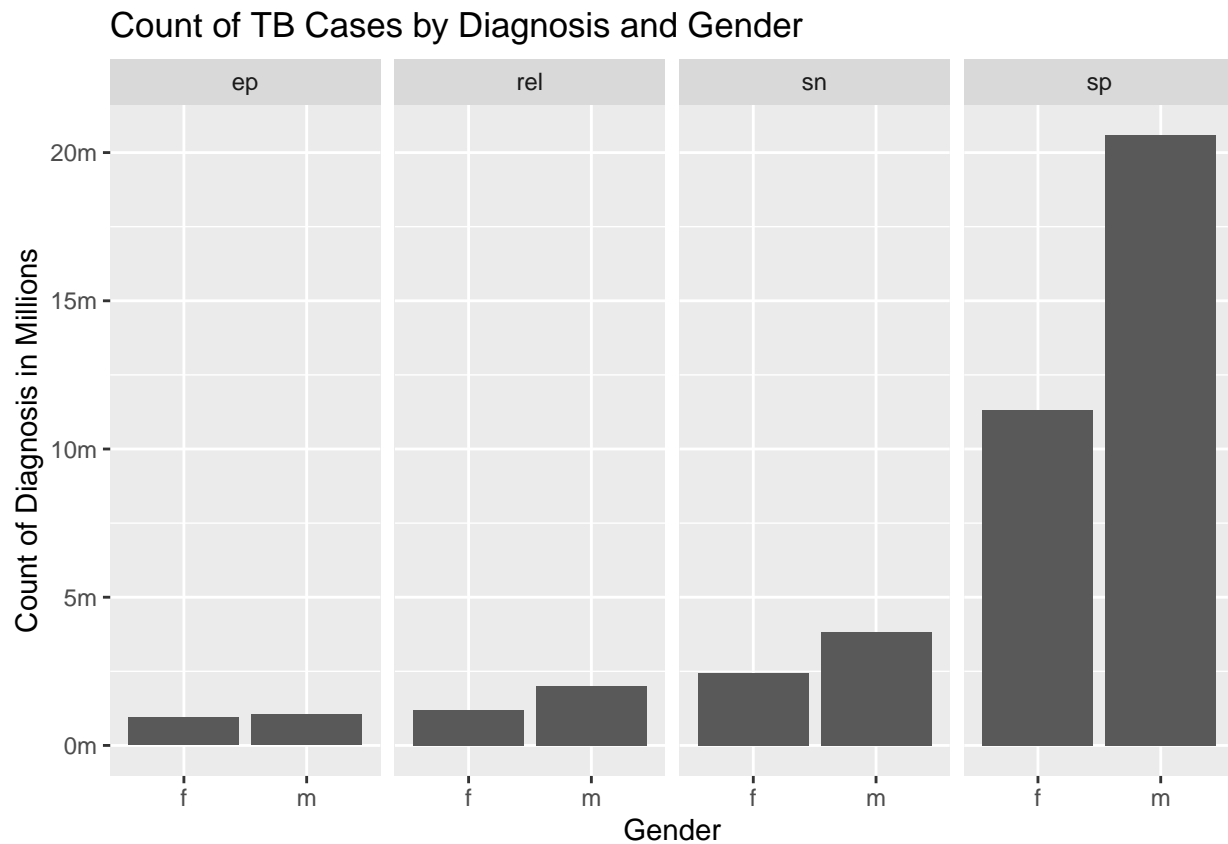
```
gender_count
```

```
## # A tibble: 8 x 3
## # Groups:   gender [2]
##   gender diagnosis count_per_diagnosis
##   <chr>   <chr>          <dbl>
## 1 f      ep              941880
## 2 f      rel             1201596
## 3 f      sn              2439139
## 4 f      sp             11324409
## 5 m      ep              1044299
## 6 m      rel             2018976
## 7 m      sn              3840388
## 8 m      sp             20586831
```

- 6) Now create a plot using ggplot and geom_col where your x axis is gender, your y axis represents the counts, and facet by diagnosis. Be sure to give your plot a title and resolve the axis labels.

```
#your code here
ggplot(gender_count) +
  geom_col(aes(x=gender, y=count_per_diagnosis)) +
```

```
facet_grid(~diagnosis) +
scale_y_continuous(labels = scales::number_format(scale=.000001, suffix="m")) +
labs(x="Gender",
     y="Count of Diagnosis in Millions",
     title = "Count of TB Cases by Diagnosis and Gender")
```



- 7) Find the percentage of population by year, gender, and diagnosis. Be sure to remove rows containing NA values.

```
# Find max population in dataset
max_pop_year <- age_split_who %>%
  group_by(year, gender, diagnosis) %>%
  #drop_na() to remove na rows
  drop_na() %>%
  #reframe instead of summarize to resolve errors (more than one result row
  #per result in join)
  reframe(max_pop=max(population, na.rm = TRUE))

# Find min population in dataset
min_pop_year <- age_split_who %>%
  group_by(year, gender, diagnosis) %>%
  drop_na() %>%
  reframe(min_pop=min(population, na.rm = TRUE))

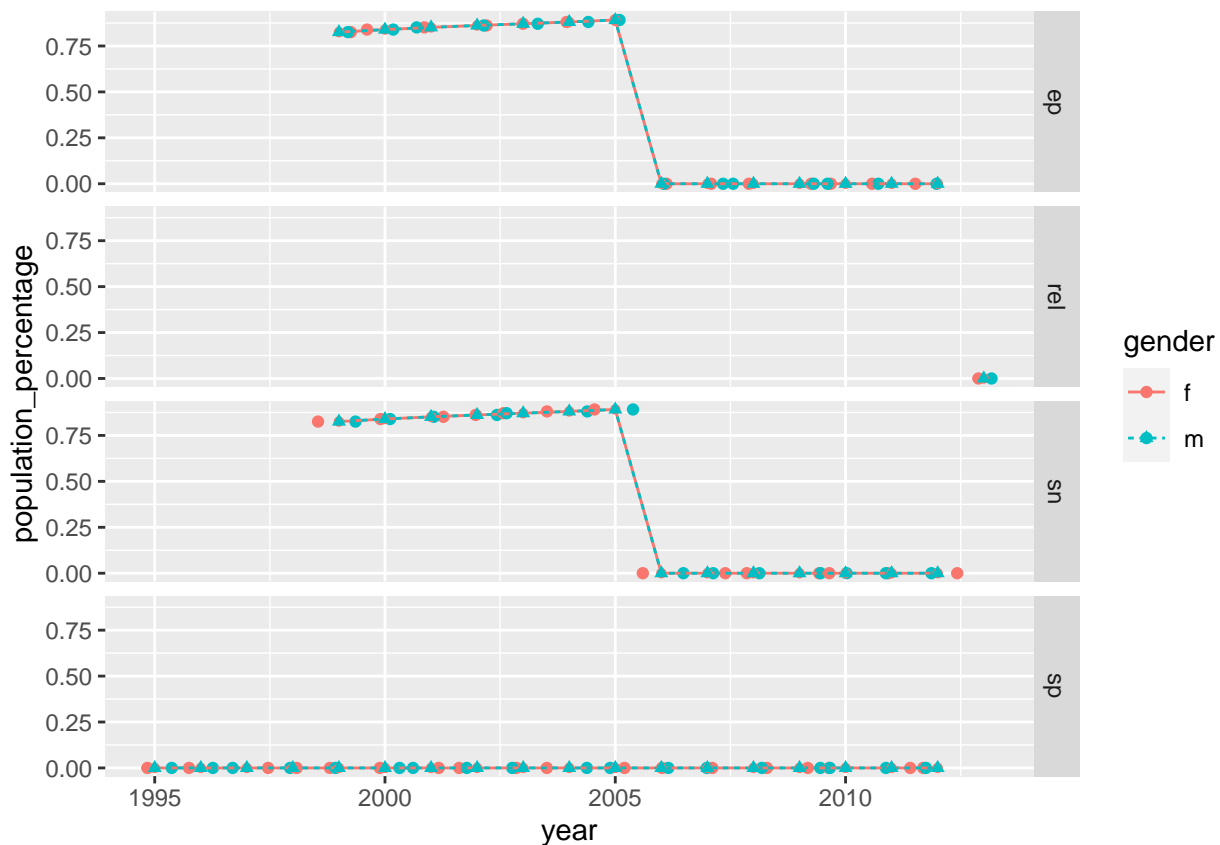
#join min and max together, and divide min over max to get percentage
age_split_who_popgroup <- max_pop_year %>%
```

```
inner_join(min_pop_year, by=c('year', 'gender', 'diagnosis')) %>%
group_by(year,gender,diagnosis) %>%
reframe(population_percentage=min_pop/max_pop)
```

- 8) Create a line plot in ggplot where your x axis contains the year and y axis contains the percent of world population. Facet this plot by diagnosis with each plot stacked vertically. You should have a line for each gender within each facet. Be sure to format your y axis and give your plot a title.

```
ggplot(age_split_who_popgroup, aes(x=year, y=population_percentage, group=gender))+
  geom_jitter(aes(color=gender), width=.5)+
  geom_line(aes(linetype=gender, color=gender))+
  geom_point(aes(shape=gender, color=gender))+
  facet_grid(rows=vars(diagnosis))
```

```
## `geom_line()`: Each group consists of only one observation.
## i Do you need to adjust the group aesthetic?
```



- 9) Now unite the min and max age variables into a new variable named age_range. Use a '-' as the separator.

```
age_split_who <- age_split_who %>%
  unite(col = 'age_range', min_age:max_age, sep='_')
head(age_split_who)
```

```
## # A tibble: 6 x 9
##   country    iso2 iso3  year diagnosis gender age_range count population
##   <chr>      <chr> <chr> <dbl> <chr>    <chr> <chr>    <dbl>    <dbl>
```

## 1	Afghanistan	AF	AFG	1980	sp	m	0_14	NA	NA
## 2	Afghanistan	AF	AFG	1980	sp	m	15_24	NA	NA
## 3	Afghanistan	AF	AFG	1980	sp	m	25_34	NA	NA
## 4	Afghanistan	AF	AFG	1980	sp	m	35_44	NA	NA
## 5	Afghanistan	AF	AFG	1980	sp	m	45_54	NA	NA
## 6	Afghanistan	AF	AFG	1980	sp	m	55_64	NA	NA

- 10) Find the percentage contribution of each age group by diagnosis. You will first need to find the count of all diagnoses then find the count of all diagnoses by age group. Join the former to the later and calculate the percent of each age group. Plot these as a `geom_col` where the x axis is the diagnosis, y axis is the percent of total, and faceted by age group.

```
#1) Find count of all diagnoses
count_diagnosis_all <- age_split_who %>%
  group_by(diagnosis) %>%
  drop_na() %>%
  summarise(total_sum=sum(count))
```

```
#2) Count of all diagnoses by age group
count_diagnosis_age <- age_split_who %>%
  group_by(diagnosis, age_range) %>%
  drop_na() %>%
  summarise(sum_by_age=sum(count))
```

``summarise()`` has grouped output by 'diagnosis'. You can override using the ## ``.groups`` argument.

```
#3) Join 1) to 2). Use inner join to find matches based on diagnosis.
count_diagnosis_combined <- count_diagnosis_all %>%
  inner_join(count_diagnosis_age, by='diagnosis') %>%
  mutate(perc_diag_by_age = (sum_by_age/total_sum))
```

```
#4) Plot data as a geom_col where the x axis is the diagnosis, y axis is the percent of total by age, a
ggplot(count_diagnosis_combined, aes(x=diagnosis,
                                     y=perc_diag_by_age))+
  geom_col(aes(fill=diagnosis))+
  facet_grid(. ~ age_range)+
  labs(x='Diagnoses',
       y='Percent By Age',
       title='Percent of Total Diagnoses by Age Group')+
  scale_y_continuous(labels = scales::percent_format(scale=100, suffix="%"))
```

