

Week 4 Exercises

Peter-Jon Grant

Nov 17th, 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>

#Load Data and Packages

```
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)
library(stringr)

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

#Pivot the data to Long format
who_longer = who |>
  pivot_longer(
    cols = starts_with("new"),
    names_to = c("diagnosis", "gender", "age"),
    names_pattern = ("new_?(.*)_(.)(.*)"),
    values_to = "count",
  )
```

- 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.

```
combined_df = who_longer |>
  left_join(
    population,
    by = c("country", "year"))

#creating a new dataframe that is easier to type
cdf = combined_df
```

- 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.

```
cdf = cdf |>
  mutate (
    age = str_replace(age, "0", "00")
  )

cdf = cdf |>
  separate(
    age,
    into = c('min_age', 'max_age'),
    sep = 2
  )
```

- 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.

```
cdf = cdf |>
  mutate(
    max_age = ifelse(max_age == "", Inf, max_age)
  )
```

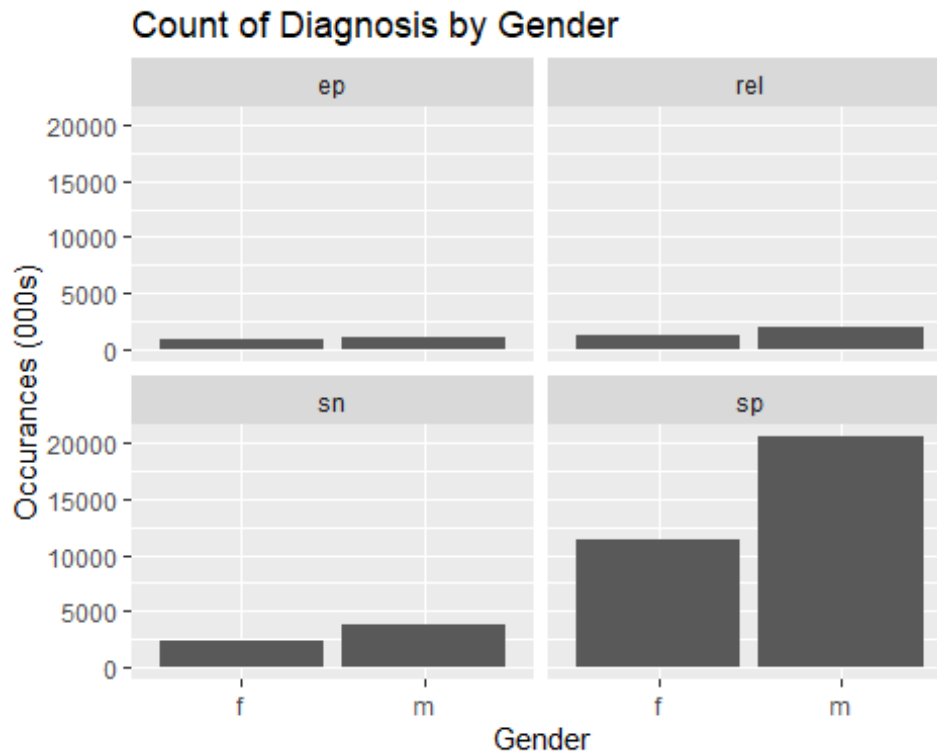
5) Find the count per diagnosis for males and females.

See ?sum for a hint on resolving NA values.

```
summary_table = cdf |>
  group_by(diagnosis, gender) |>
  summarise(
    count = sum(count, na.rm = TRUE) / 1000,
    .groups = 'drop'
  )
```

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.

```
ggplot(summary_table, aes(x = gender, y = count)) +
  geom_col() +
  facet_wrap(~diagnosis) +
  labs(
    title = "Count of Diagnosis by Gender",
    x = 'Gender',
    y = 'Occurances (000s)'
  )
```



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

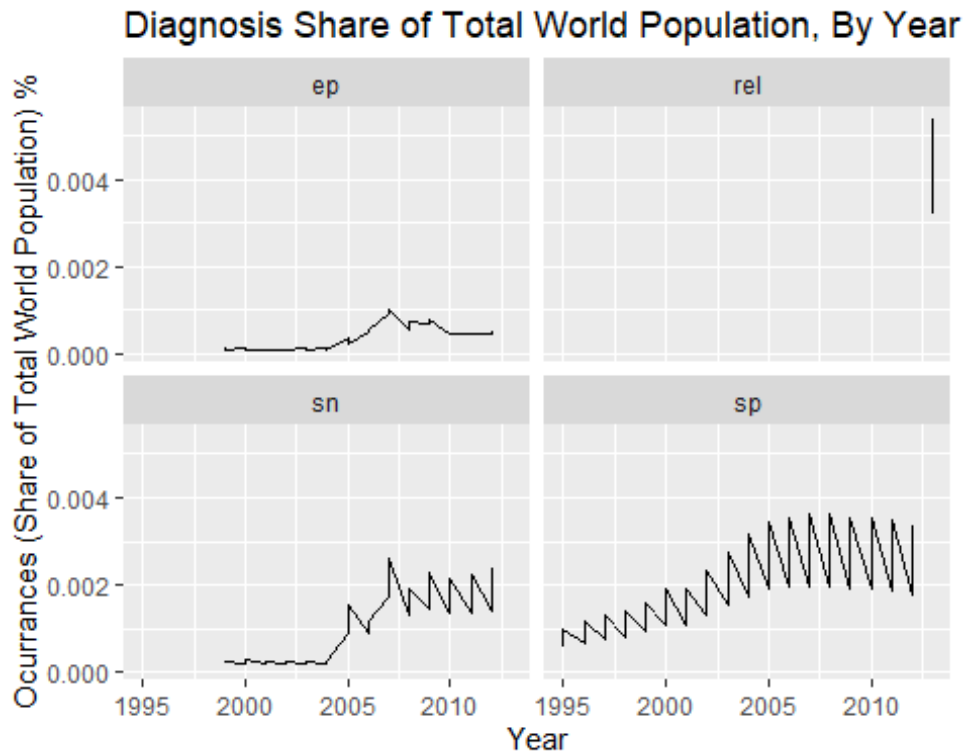
```
#Remove rows with NAs
```

```
cdf = na.omit(cdf)
```

```
population_stats = cdf |>
  group_by(year, gender, diagnosis)|>
  summarise(
    world_pop = sum(population),
    occurrences = sum(count),
    share_of_pop = (occurrences / world_pop) * 100,
    .groups = 'drop'
  )
```

- 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(population_stats, aes(year, share_of_pop)) + geom_line() +
  facet_wrap(~diagnosis) + labs(
    title = "Diagnosis Share of Total World Population, By Year",
    x = 'Year',
    y = 'Occurrences (Share of Total World Population) %'
  )
```



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

```
cdf$age_range = paste(cdf$min_age, cdf$max_age, sep = "-")
```

- 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.

```
age_share = cdf

age_share = age_share |>
  group_by(diagnosis) |>
  mutate(
    total_diagnosis_count = sum(count)
  ) |>
  group_by(age_range, diagnosis) |>
  mutate(
    age_count_diagnosis = sum(count),
    age_share_diagnosis = (age_count_diagnosis / total_diagnosis_count) *
100
  )

ggplot(age_share, aes(x = diagnosis, y = age_share_diagnosis, fill =
age_range)) +
  geom_col(position = "dodge") +
```

```
labs(
  title = "Age Share by Diagnosis",
  x = "Diagnosis",
  y = "Age Share of Diagnosis (%)"
) +
facet_wrap(~ age_range)
```

