

STA2201_Lab1

Alice Huang

11/01/2023

By the end of this lab you should know the basics of

- RStudio Projects
- Quarto
- Main tidyverse functions
- ggplot

RStudio Projects

RStudio projects are associated with R working directories. They are good to use for several reasons:

- Each project has their own working directory, so make dealing with file paths easier
- Make it easy to divide your work into multiple contexts
- Can open multiple projects at one time in separate windows

To make a new project in RStudio, go to File → New Project. If you've already set up a repo for this class, then select 'Existing Directory' and choose the folder that will contain all your class materials. This will open a new RStudio window, that will be called the name of your folder.

In future, when you want to do work for this class, go to your class folder and open the .Rproj file. This will open an RStudio window, with the correct working directory, and show the files you were last working on.

Quarto

This is a Quarto document. Quarto allows you to create nicely formatted documents (HTML, PDF, Word) that also include code and output of the code. This is good because it's reproducible, and also makes reports easier to update when new data comes in. Each of the grey chunks contains R code, just like a normal R script. You can choose to run each chunk separately, or knit the whole document using Knit the button above, which creates your document.

To start a new Quarto file in Rstudio, go to File → New File → Quarto, then select Document and whatever you want to compile the document as (I chose pdf, and that's generally what we'll be doing in this class). Notice that this and the other inputs (title, author) are used to create the 'yaml', the bit at the start of the document. You can edit this, like I have for example to include table of contents.

When you hit 'Render' a pdf will be created and saved in the same folder as your `qmd` file. There are various options for output code, results, etc. For example, if you don't want your final report to include the code (but just the output, e.g. graphs or tables) then you can specify `#| echo=FALSE` at the top of the chunk (note: this syntax is different R Markdown).

Quarto is a newer version of R Markdown. It's the first year I'm switching and so I'm fine if people would rather stick to R Markdown for now. A helpful intro is here: <https://quarto.org/docs/get-started/hello/rstudio.html>

Writing math

Writing equations is essentially the same as in LaTeX. You can write inline equations using the $\$$ e.g. $y = ax + b$. You can write equations on a separate line with two $\$$ s e.g.

$$y = ax + b$$

In pdf documents you can have numbered equations using

$$y = ax + b \tag{1}$$

Getting greek letters, symbols, subscripts, bars etc is the same as LaTeX. A few examples are below

- $Y_{i,j}$
- $\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$
- $\alpha\beta\gamma$
- $X \rightarrow Y$
- $Y \sim N(\mu, \sigma^2)$

Tidyverse

Read in some packages that we'll be using:

```
#install.packages("tidyverse")
library(tidyverse)
```

On top of the base R functionality, there's lots of different packages that different people have made to improve the usability of the language. One of the most successful suite of packages is now called the 'tidyverse'. The tidyverse contains a range of functionality that help to manipulate and visualize data.

Read in mortality rates for Ontario. These data come from the Canadian Human Mortality Database.

```
dm <- read_table("https://www.prhdh.umontreal.ca/BDLC/data/ont/Mx_1x1.txt", skip = 2, col_types = "dcddd")
head(dm)
```

```
## # A tibble: 6 x 5
##   Year Age   Female   Male   Total
##   <dbl> <chr>   <dbl>   <dbl>   <dbl>
## 1  1921 0     0.0978  0.129   0.114
## 2  1921 1     0.0129  0.0144  0.0137
## 3  1921 2     0.00521 0.00737 0.00631
## 4  1921 3     0.00471 0.00457 0.00464
## 5  1921 4     0.00461 0.00433 0.00447
## 6  1921 5     0.00372 0.00361 0.00367
```

The object `dm` is a data frame, or tibble. Every column can be a different data type (e.g. we have integers and characters).

Important tidyverse functions

You should feel comfortable using the following functions

- The pipe `|>` or `%>%`
- `filter`
- `select`
- `arrange`
- `mutate`
- `group_by`
- `summarize`
- `pivot_longer` and `pivot_wider`

Piping, filtering, selecting, arranging

A central part of manipulating tibbles is using the `|>` function. This is a pipe, but should be read as saying ‘and then’. Note that the `|>` syntax is the base version of the pipe (new last year). Previously the syntax was `%>%` through the `magrittr` package. They essentially function the same.

For example, say we just want to pull out mortality rates for 1935. We would take our tibble *and then* filter to only include 1935:

```
dm %>%  
  filter(Year==1935) # two equals signs logical
```

```
## # A tibble: 111 x 5  
##   Year Age   Female   Male   Total  
##   <dbl> <chr>   <dbl>   <dbl>   <dbl>  
## 1  1935 0     0.0513  0.0652  0.0584  
## 2  1935 1     0.00607 0.00742 0.00676  
## 3  1935 2     0.00350 0.00321 0.00336  
## 4  1935 3     0.00187 0.00321 0.00255  
## 5  1935 4     0.0013   0.00238 0.00185  
## 6  1935 5     0.00152 0.00186 0.00169  
## 7  1935 6     0.00136 0.00174 0.00155  
## 8  1935 7     0.00120 0.00154 0.00137  
## 9  1935 8     0.000984 0.00130 0.00114  
## 10 1935 9     0.000996 0.00140 0.00120  
## # ... with 101 more rows
```

You can also filter by more than one condition; say we just wanted to look at 10 year olds in 1935:

```
dm %>% filter(Year ==1935, Age == 10)
```

```
## # A tibble: 1 x 5  
##   Year Age   Female   Male   Total  
##   <dbl> <chr>   <dbl>   <dbl>   <dbl>  
## 1  1935 10     0.000884 0.00143 0.00116
```

If we only wanted to look at 10 year olds in 1935 who were female, we could filter *and then* select the female column.

```
dm %>% filter(Year == 1935, Age == 10,) %>% select(Female)
```

```
## # A tibble: 1 x 1
##   Female
##   <dbl>
## 1 0.000884
```

You can also remove columns by selecting the negative of that column name.

```
dm %>% select(-Female)
```

```
## # A tibble: 10,989 x 4
##   Year Age      Male      Total
##   <dbl> <chr>   <dbl>   <dbl>
## 1 1921 0      0.129   0.114
## 2 1921 1      0.0144  0.0137
## 3 1921 2      0.00737 0.00631
## 4 1921 3      0.00457 0.00464
## 5 1921 4      0.00433 0.00447
## 6 1921 5      0.00361 0.00367
## 7 1921 6      0.00393 0.00330
## 8 1921 7      0.00351 0.00323
## 9 1921 8      0.00285 0.00262
## 10 1921 9      0.00255 0.00227
## # ... with 10,979 more rows
```

Sort the tibble according to a particular column using `arrange`, for example, Year in descending order:

```
dm %>% arrange(desc(Year))
```

```
## # A tibble: 10,989 x 5
##   Year Age      Female      Male      Total
##   <dbl> <chr>   <dbl>   <dbl>   <dbl>
## 1 2019 0      0.00423 0.00481 0.00453
## 2 2019 1      0.000216 0.000177 0.000196
## 3 2019 2      0.000157 0.000162 0.00016
## 4 2019 3      0.00007 0.00016 0.000117
## 5 2019 4      0.000111 0.000132 0.000122
## 6 2019 5      0.000096 0.000052 0.000074
## 7 2019 6      0.000081 0.000039 0.000059
## 8 2019 7      0.000107 0.000128 0.000118
## 9 2019 8      0.000066 0.000026 0.000046
## 10 2019 9      0.000052 0.000177 0.000116
## # ... with 10,979 more rows
```

NOTE: none of the above operations are saving. To save, you need to assign the output to an object. You can call it something new or overwrite the original.

```
dm1935 <- dm %>% filter(Year==1935)
```

Grouping, summarizing, mutating

In addition to `filter` and `select`, two useful functions are `mutate`, which allows you to create new variables, and `summarize`, which allows you to produce summary statistics. These are particularly powerful when combined with `group_by()` which allows you to do any operation on a tibble by group.

For example, let's create a new variable that is the ratio of male to female mortality at each age and year:

```
dm %>%  
  mutate(mf_ratio = Male/Female)
```

```
## # A tibble: 10,989 x 6  
##   Year Age   Female   Male   Total mf_ratio  
##   <dbl> <chr>   <dbl>   <dbl>   <dbl>   <dbl>  
## 1 1921 0     0.0978  0.129   0.114    1.32  
## 2 1921 1     0.0129  0.0144  0.0137    1.11  
## 3 1921 2     0.00521 0.00737 0.00631    1.41  
## 4 1921 3     0.00471 0.00457 0.00464    0.970  
## 5 1921 4     0.00461 0.00433 0.00447    0.938  
## 6 1921 5     0.00372 0.00361 0.00367    0.971  
## 7 1921 6     0.00265 0.00393 0.00330    1.48  
## 8 1921 7     0.00295 0.00351 0.00323    1.19  
## 9 1921 8     0.00237 0.00285 0.00262    1.20  
## 10 1921 9     0.00198 0.00255 0.00227    1.29  
## # ... with 10,979 more rows
```

Now, let's calculate the mean female mortality rate by age over all the years. To do this, we need to `group_by` Age, and then use `summarize` to calculate the mean:

```
dm %>% group_by(Age) %>%  
  summarize(mean = mean(Female))
```

```
## # A tibble: 111 x 2  
##   Age   mean  
##   <chr> <dbl>  
## 1 0      NA  
## 2 1      NA  
## 3 10     NA  
## 4 100    NA  
## 5 101    NA  
## 6 102    NA  
## 7 103    NA  
## 8 104    NA  
## 9 105    NA  
## 10 106   NA  
## # ... with 101 more rows
```

Mean of males and females by age

```
dm %>% group_by(Age) %>%  
  summarize(mean_female = mean(Female, na.rm=TRUE), mean_male = mean(Male, na.rm=TRUE))
```

```
## # A tibble: 111 x 3
##   Age    mean_female mean_male
##   <chr>      <dbl>      <dbl>
## 1 0          0.0254      0.0322
## 2 1          0.00262     0.00297
## 3 10         0.000426     0.000590
## 4 100        0.426       0.462
## 5 101        0.448       0.493
## 6 102        0.493       0.566
## 7 103        0.533       0.647
## 8 104        0.660       0.780
## 9 105        0.805       0.904
## 10 106       0.796       0.720
## # ... with 101 more rows
```

Alternatively using across

```
dm %>% group_by(Age) %>%
  summarize(across(c("Male", "Female"), mean))
```

```
## # A tibble: 111 x 3
##   Age      Male    Female
##   <chr>    <dbl>    <dbl>
## 1 0      0.0322    0.0254
## 2 1      0.00297    0.00262
## 3 10     0.000590    0.000426
## 4 100    0.462     0.426
## 5 101    0.493     0.448
## 6 102    0.566     0.493
## 7 103    0.647     0.533
## 8 104    NA        0.660
## 9 105    NA        NA
## 10 106    NA        NA
## # ... with 101 more rows
```

Pivoting

We often need to switch between wide and long data format. The `dm` tibble is currently in wide format. To get it in long format we can use `pivot_longer`

```
dm_long <- dm %>%
  pivot_longer(Female:Total, names_to = "sex", values_to = "mortality")
dm_long
```

```
## # A tibble: 32,967 x 4
##   Year Age sex    mortality
##   <dbl> <chr> <chr>      <dbl>
## 1 1921 0   Female  0.0978
## 2 1921 0   Male    0.129
## 3 1921 0   Total   0.114
## 4 1921 1   Female  0.0129
## 5 1921 1   Male    0.0144
```

```
## 6 1921 1      Total    0.0137
## 7 1921 2      Female   0.00521
## 8 1921 2      Male     0.00737
## 9 1921 2      Total    0.00631
## 10 1921 3     Female   0.00471
## # ... with 32,957 more rows
```

```
dm_long %>% pivot_wider(names_from = sex, values_from = "mortality") -> dm
dm
```

```
## # A tibble: 10,989 x 5
##   Year Age      Female      Male      Total
##   <dbl> <chr>    <dbl>    <dbl>    <dbl>
## 1 1921 0      0.0978  0.129    0.114
## 2 1921 1      0.0129  0.0144   0.0137
## 3 1921 2      0.00521 0.00737  0.00631
## 4 1921 3      0.00471 0.00457  0.00464
## 5 1921 4      0.00461 0.00433  0.00447
## 6 1921 5      0.00372 0.00361  0.00367
## 7 1921 6      0.00265 0.00393  0.00330
## 8 1921 7      0.00295 0.00351  0.00323
## 9 1921 8      0.00237 0.00285  0.00262
## 10 1921 9     0.00198 0.00255  0.00227
## # ... with 10,979 more rows
```

Using ggplot

You can plot things in R using the base `plot` function, but plots using `ggplot` are much prettier.

Say we wanted to plot the mortality rates for 30 year old males over time. In the function `ggplot`, we need to specify our data (in this case, a filtered version of `dm`), an x axis (Year) and y axis (Male). The axes are defined withing the `aes()` function, which stands for 'aesthetics'.

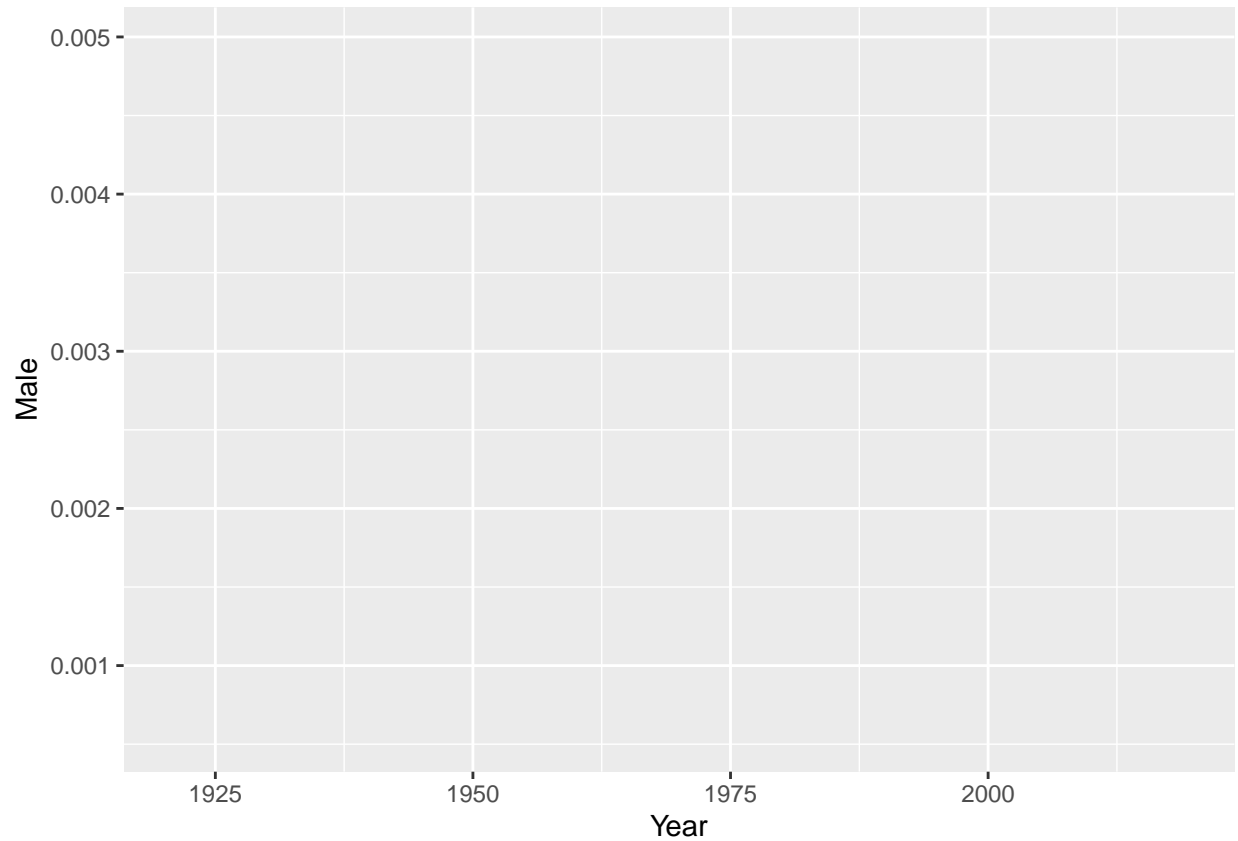
First let's get our data:

```
d_to_plot <- dm %>%
  filter(Age==30) %>%
  select(Year, Male)
d_to_plot
```

```
## # A tibble: 99 x 2
##   Year      Male
##   <dbl>    <dbl>
## 1 1921 0.00375
## 2 1922 0.00462
## 3 1923 0.00497
## 4 1924 0.00412
## 5 1925 0.00308
## 6 1926 0.00308
## 7 1927 0.00327
## 8 1928 0.00356
## 9 1929 0.00393
## 10 1930 0.00418
## # ... with 89 more rows
```

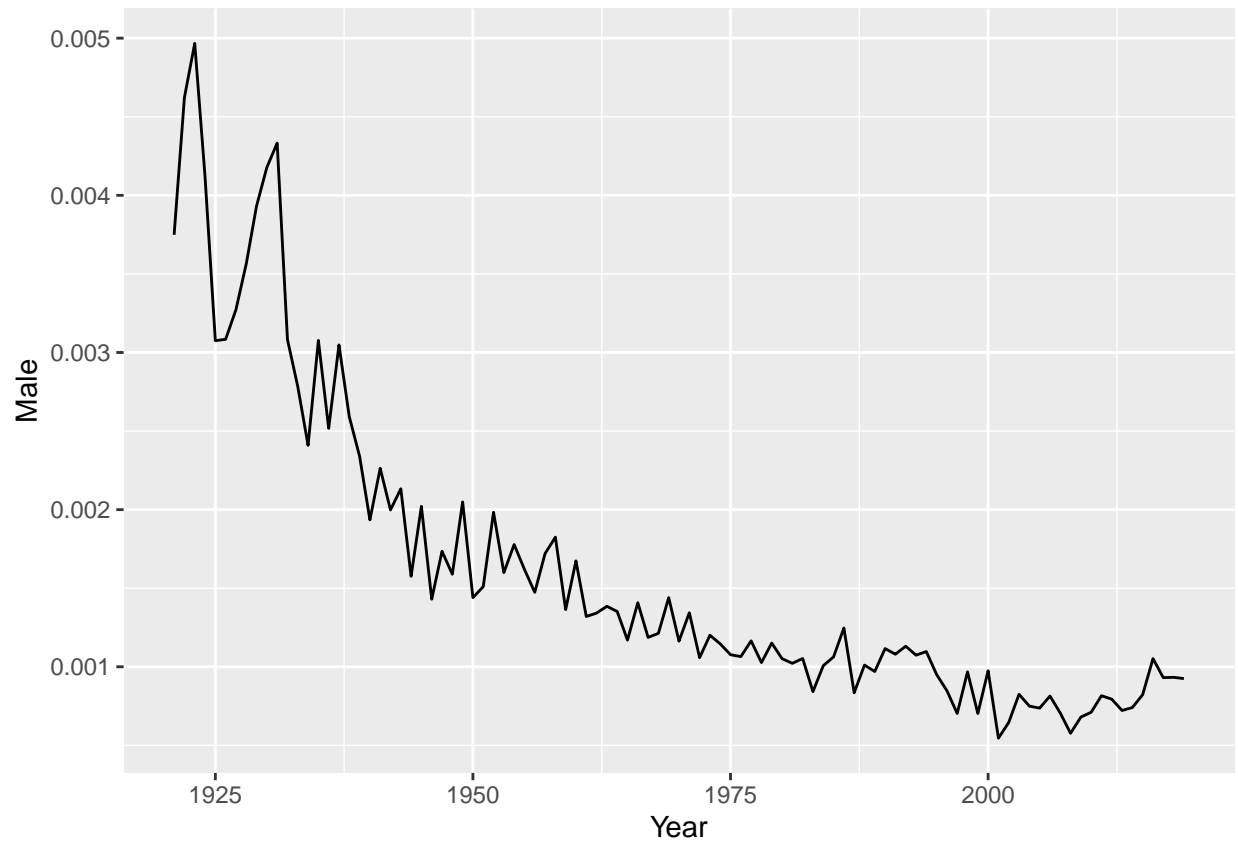
Now start the ggplot:

```
p <- ggplot(data = d_to_plot, aes(x = Year, y = Male))  
p
```



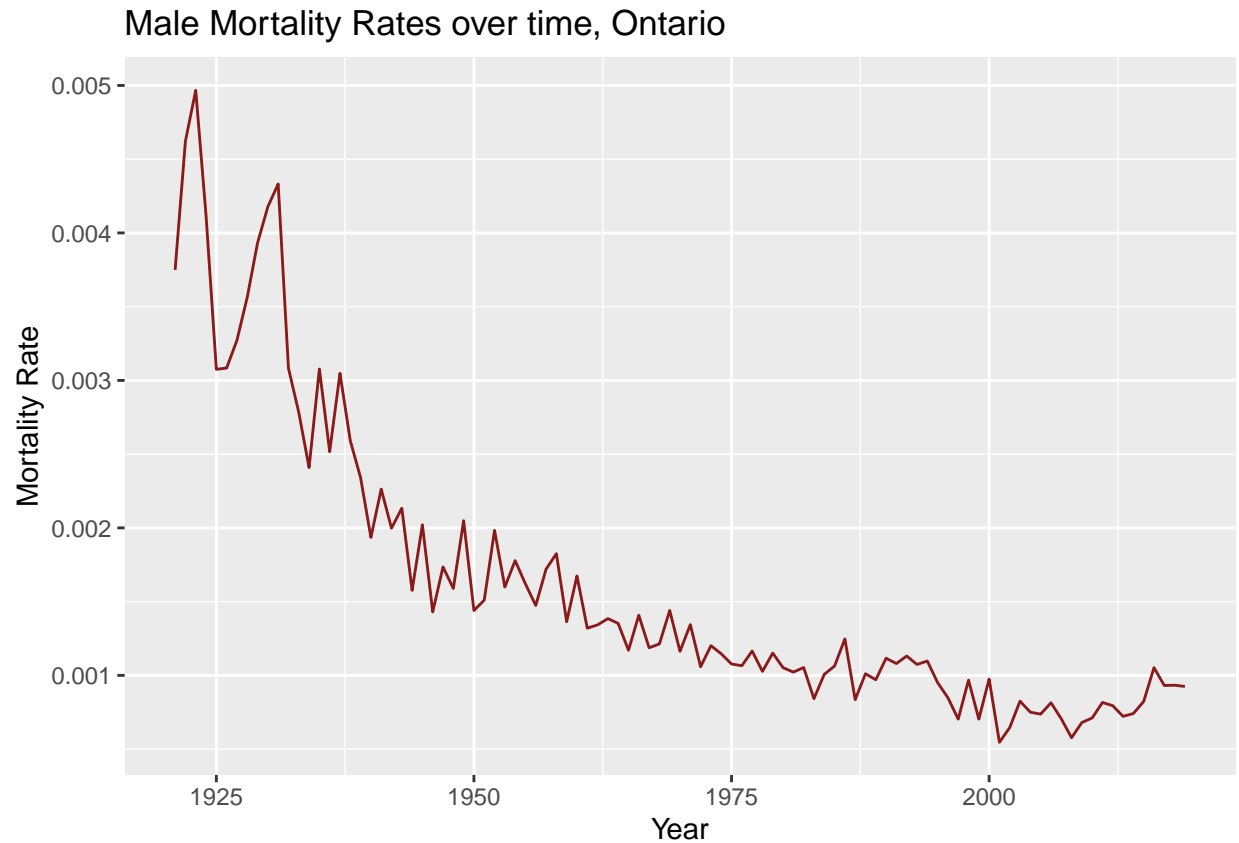
Notice the object `p` is just an empty box. The key to ggplot is layering: we now want to specify that we want a line plot using `geom_line()`:

```
p <- ggplot(data = d_to_plot, aes(x = Year, y = Male)) + geom_line()  
p
```

Let's change the color of the line, and the y-axis label, and give the plot a title:

```
p <- ggplot(data = d_to_plot, aes(x = Year, y = Male)) + geom_line(color = "firebrick4") + labs(title =  
p
```



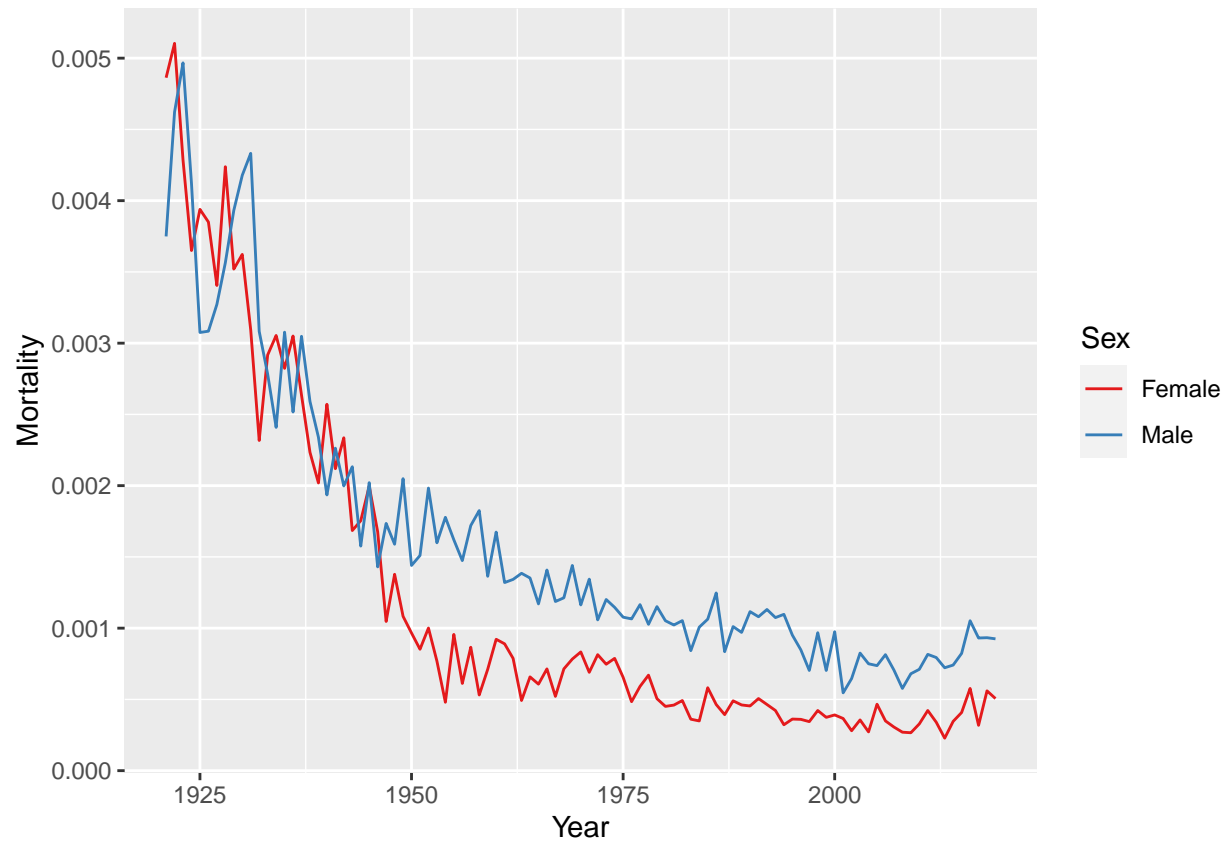
More than one group

Now say we wanted to have trends for 30-year old males and females on the one plot. The easiest way to do this is to first reshape our data so it's in long format: so instead of having a column for each sex, we have one column indicating the sex, and another column indicating the Mx value

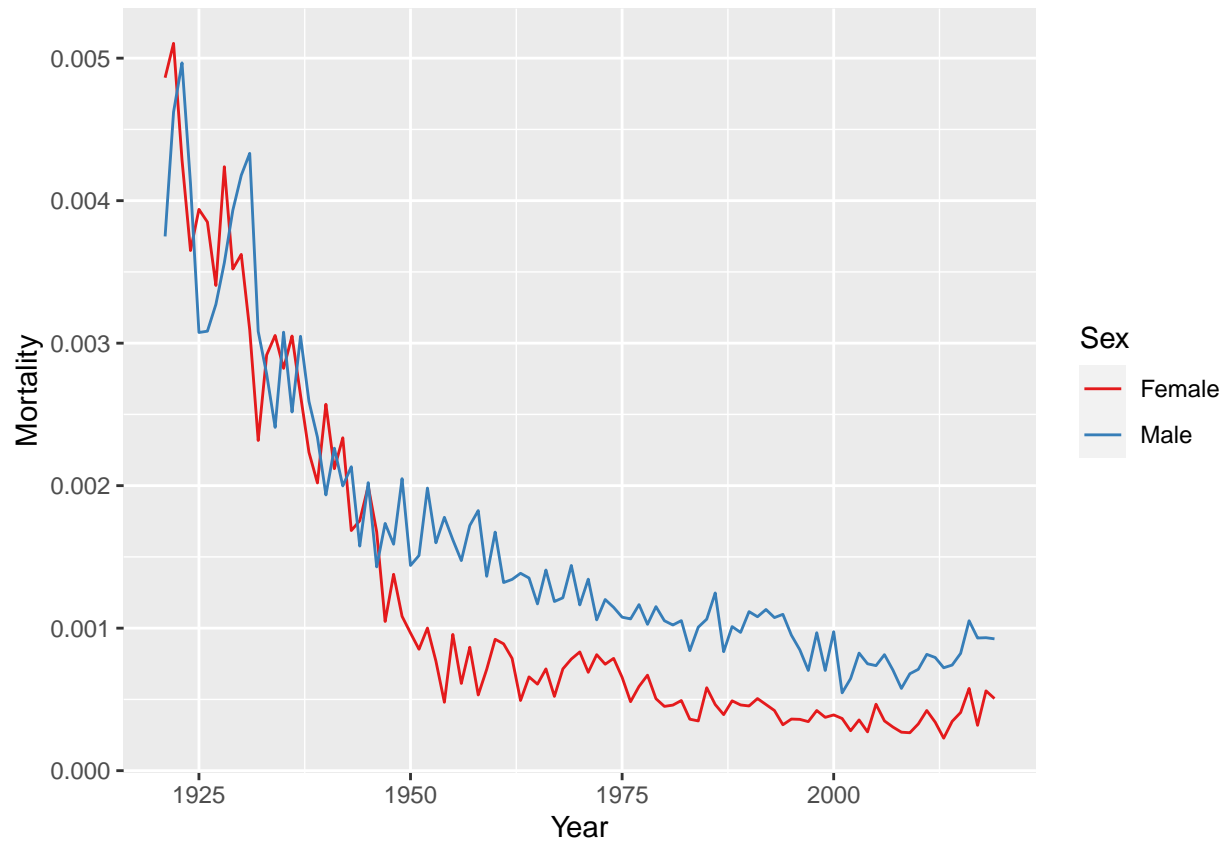
```
dm %>% filter(Age == 30) %>%
  select(Year:Male) %>%
  pivot_longer(Female:Male, names_to = "Sex", values_to = "Mortality") -> dp
```

Now we can do a similar plot to before but we now have an added component in the `aes()` function: `color`, which is determined by `sex`:

```
dp %>% ggplot(aes(x=Year, y=Mortality, color = Sex)) + geom_line() + scale_color_brewer(palette="Set1")
```



```
dp %>% ggplot(aes(x=Year, y=Mortality, color = Sex)) + geom_line() + scale_color_brewer(palette="Set1")
```



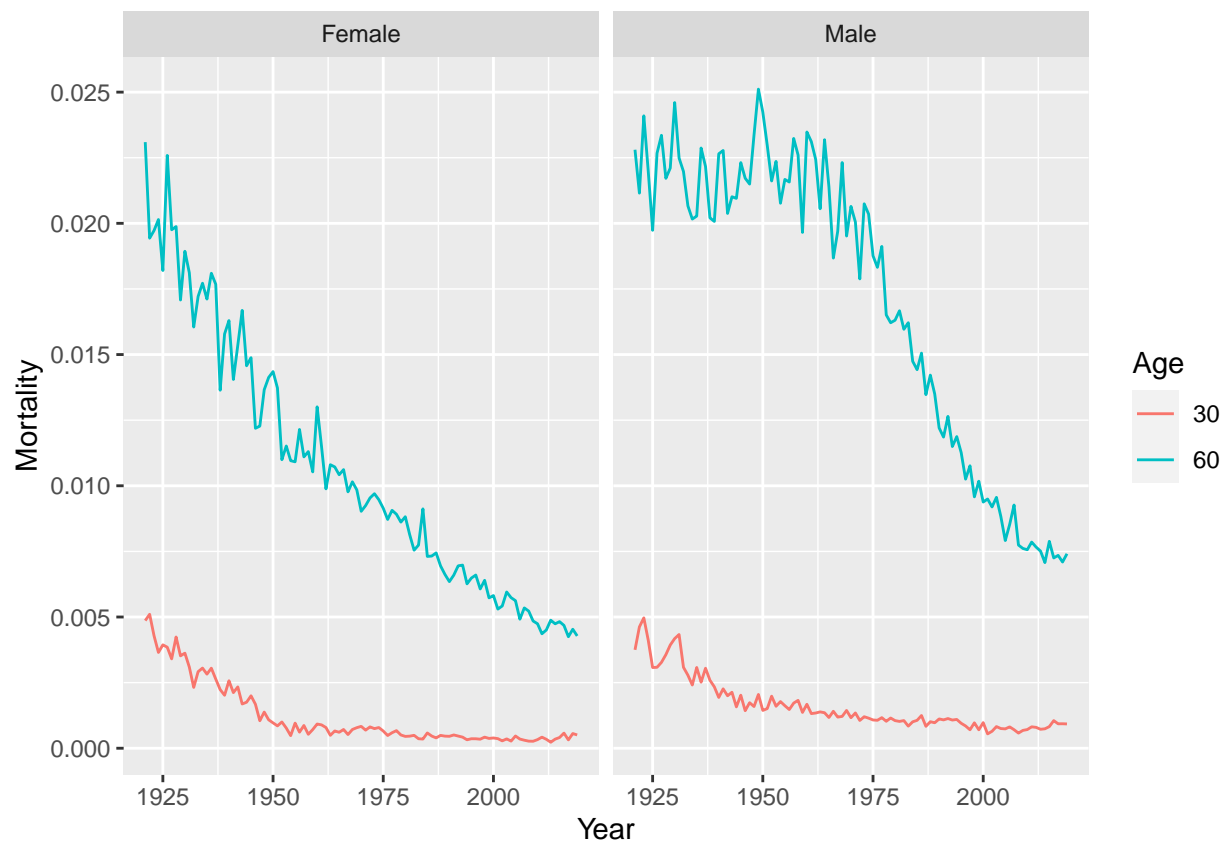
Faceting

A neat thing about ggplot is that it's relatively easy to create 'facets' or smaller graphs divided by groups. Say we wanted to look at trends for 30 year olds and 60 year olds for both males and females. Let's get the data ready to plot:

```
dm %>% filter(Age == 30 | Age == 60) %>%
  select(Year:Male) %>%
  pivot_longer(Female:Male, names_to = "Sex", values_to = "Mortality") -> dp
```

Now let's plot, with a separate facet for each sex:

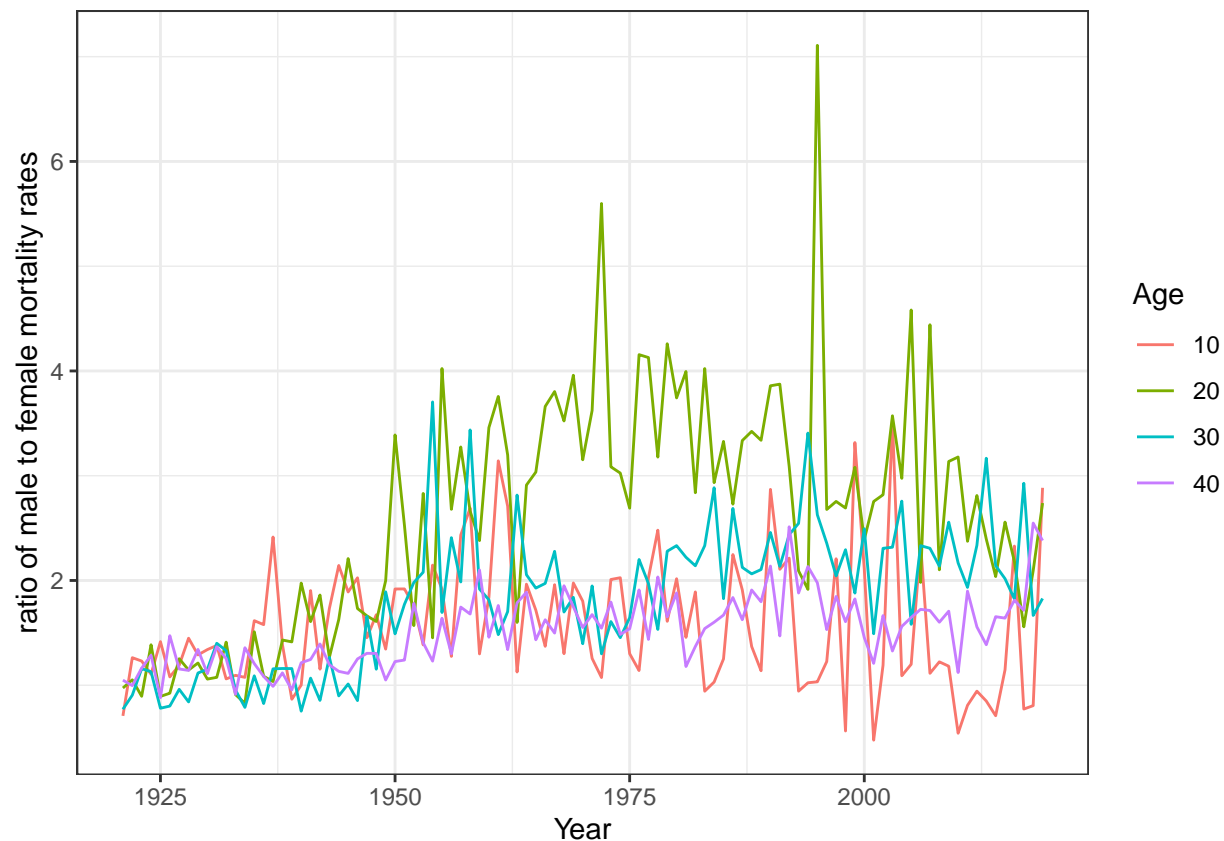
```
dp %>% ggplot(aes(x=Year, y=Mortality, color = Age)) + geom_line() + facet_grid(~Sex)
```



Lab Exercises

1. Plot the ratio of male to female mortality rates over time for ages 10,20,30 and 40 (different color for each age) and change the theme

```
dm %>% filter(Age == 10 | Age == 20 | Age == 30 | Age == 40)%>%
  mutate(mf_ratio = Male/Female) %>%
  ggplot(aes(x=Year, y= mf_ratio, col=Age)) + geom_line() + theme_bw() + labs(y="ratio of male to female")
```



2. Find the age that has the highest female mortality rate each year

```
dm %>% select(Year, Age, Female) %>%
  group_by(Year) %>%
  # subtract 1 since age starts at 0, and R start indexing at 1
  summarize(m = which.max(Female)-1)
```

```
## # A tibble: 99 x 2
##   Year      m
##   <dbl> <dbl>
## 1 1921   106
## 2 1922    98
## 3 1923   104
## 4 1924   107
## 5 1925    98
## 6 1926   106
## 7 1927   106
## 8 1928   104
## 9 1929   104
## 10 1930   105
## # ... with 89 more rows
```

3. Use the `summarize(across())` syntax to calculate the standard deviation of mortality rates by age for the Male, Female and Total populations.

```
dm %>% summarize(across(c("Male", "Female", "Total"), sd, na.rm=TRUE))
```

```
## # A tibble: 1 x 3
##   Male Female Total
##   <dbl> <dbl> <dbl>
## 1 0.333 0.311 0.296
```

The standard deviations of mortality rates for the Male, Female and Total populations are respectively 0.3327571, 0.3111401 and 0.2960847.

4. The Canadian HMD also provides population sizes over time (<https://www.prddh.umontreal.ca/BDLC/data/ont/Population.txt>). Use these to calculate the population weighted average mortality rate separately for males and females, for every year. Make a nice line plot showing the result (with meaningful labels/titles) and briefly comment on what you see (1 sentence). Hint: `left_join` will probably be useful here.

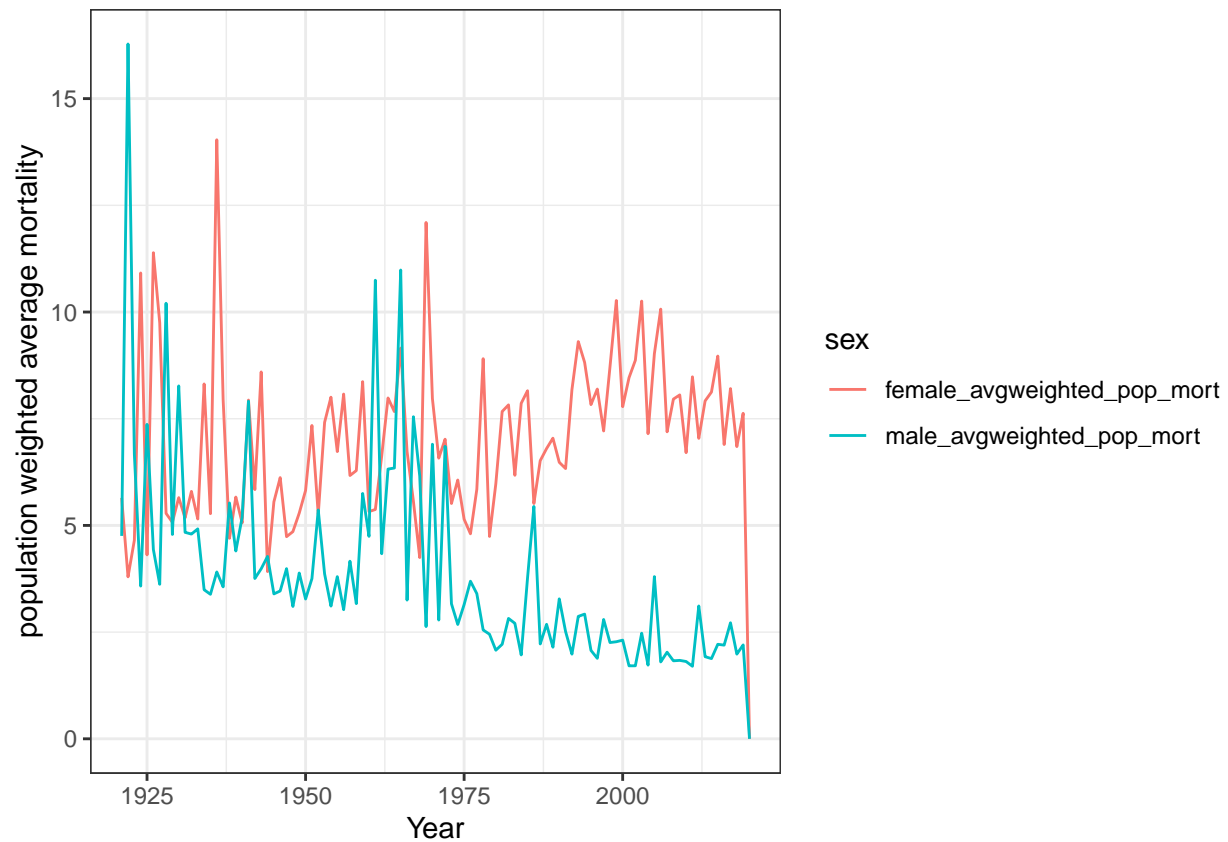
```
pop <- read_table("https://www.prddh.umontreal.ca/BDLC/data/ont/Population.txt", skip = 2, col_types = "d")
pop_dm <- left_join(pop, dm, by=c("Year", "Age"), suffix = c("_pop", "_mort"))
head(pop_dm)
```

```
## # A tibble: 6 x 8
##   Year Age   Female_pop Male_pop Total_pop Female_mort Male_mort Total_mort
##   <dbl> <chr>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 1921 0       30157.  31530.  61687.  0.0978  0.129  0.114
## 2 1921 1       30391.  31319.  61711.  0.0129  0.0144 0.0137
## 3 1921 2       30962.  31785.  62747.  0.00521 0.00737 0.00631
## 4 1921 3       31306.  32031.  63336.  0.00471 0.00457 0.00464
## 5 1921 4       31364.  32046.  63409.  0.00461 0.00433 0.00447
## 6 1921 5       31175.  31847.  63021.  0.00372 0.00361 0.00367
```

```
pop_dm %>% mutate(Female_popweight_mort = (Female_pop/Total_pop)*Female_mort) %>%
  mutate(Male_popweight_mort = (Male_pop/Total_pop)*Male_mort) %>%
  group_by(Year) %>%
  summarize(female_avgweighted_pop_mort = sum(Female_popweight_mort, na.rm=TRUE),
            male_avgweighted_pop_mort = sum(Male_popweight_mort, na.rm=TRUE)) -> newpop
head(newpop)
```

```
## # A tibble: 6 x 3
##   Year female_avgweighted_pop_mort male_avgweighted_pop_mort
##   <dbl> <dbl> <dbl>
## 1 1921 5.65 4.75
## 2 1922 3.80 16.3
## 3 1923 4.65 6.67
## 4 1924 10.9 3.58
## 5 1925 4.31 7.37
## 6 1926 11.4 4.43
```

```
newpop %>% pivot_longer(c(female_avgweighted_pop_mort, male_avgweighted_pop_mort), names_to = "sex", values_to = "mortality") %>%
  ggplot(aes(x=Year, y=mortality, color = sex)) + geom_line() + theme_bw() + labs(y="population weighted average mortality")
```



From 1975-2020, it looks like the gap between the Female and Male average population-weighted mortality is increasing. The male average population-weighted mortality seems to show a decreasing trend, but the same cannot be said for the female average population-weighted mortality.