

# Bespoke Data Visualization using R and ggplot2

A CHI 2019 Course

*Sandy Gould*

*University of Birmingham*

*s.gould@cs.bham.ac.uk*

*Glasgow, 8<sup>th</sup> May 2019*



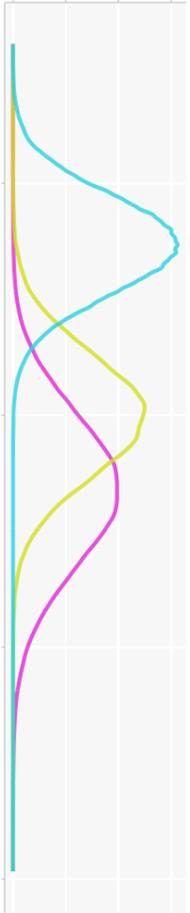
# This course

## Welcome!

## Thank you for joining this course!

There are two parts to it:

- I want to talk to you about our approach to visualization in this course.
- Then I want us to work through a set of practical tasks. This will be the vast majority of our time.



# Glossary

Term	Meaning
R	A statistical programming language
RStudio	An environment for writing R programs
ggplot2	A library for creating visualizations using R
Declarative visualization	A way of specifying visualizations



Materials

**<http://sjjg.uk/chi19-course>**

Click the link for materials. You'll need credentials.

*Username:*            **chi2019**

*Password:*            **ggplot**



# Declarative visualizations

Specifying visualizations declaratively means that:

- Output is **predictable**
  - *e.g., if we want to redraw*
- It's easy to completely change the visualization without having to mess around with the underlying data.
- We can create more complex visualizations with a higher degree of **reliability**
- Text-based declarations will always be **readable**



# The ggplot2 approach

ggplot2 takes a declarative approach to the creation of visualizations.

The emphasis is on:

- A logical structure to the creation of visualizations
- Structuring of the data that is being used
- High quality defaults



# 'The Grammar of Graphics'

Hadley Wickham, the developer of of ggplot2 has published on the philosophy of ggplot2 in his book 'The Grammar of Graphics'.

Essentially, everything is split into:

- **Aesthetics** (*what is going to be shown and how it will be shown*)
- **Geometries** (*the form of plots that will represent the data*)

*I find the labels for these confusing at times!*



# Writing code for R

- R is written using plaintext files.
- They can be written using a text editor and run from the command line.
- They can be written and run from R's default editing environment.
- They can be written and run from RStudio, a complete development environment for R.



# Writing code for R

- The goal of today is not to teach you programming.
- It is to give you an understanding of how plots can be pieced together using a variety of ggplot2 commands.
- You will leave with an understanding of the logic of ggplot2 as well as plenty of template code to get you started.
- The website will stay up for you to refer to in future.



# RStudio

The screenshot shows the RStudio interface with several panels highlighted by colored boxes:

- Yellow box:** The Environment, History, and Connections panels. The Environment panel shows "Global Environment" and "Environment is empty".
- Green box:** The Files, Plots, Packages, Help, and Viewer panels. The Files panel shows a file named "kitematic" under the "Home" directory.
- Red box:** The Console panel, which displays the R version information and license details.

```
1
```

1:1 (Top Level) R Script

Console Terminal Jobs

```
~/
```

R version 3.5.3 (2019-03-11) -- "Great Truth"  
Copyright (C) 2019 The R Foundation for Statistical Computing  
Platform: x86\_64-pc-linux-gnu (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.



# Getting RStudio up and running

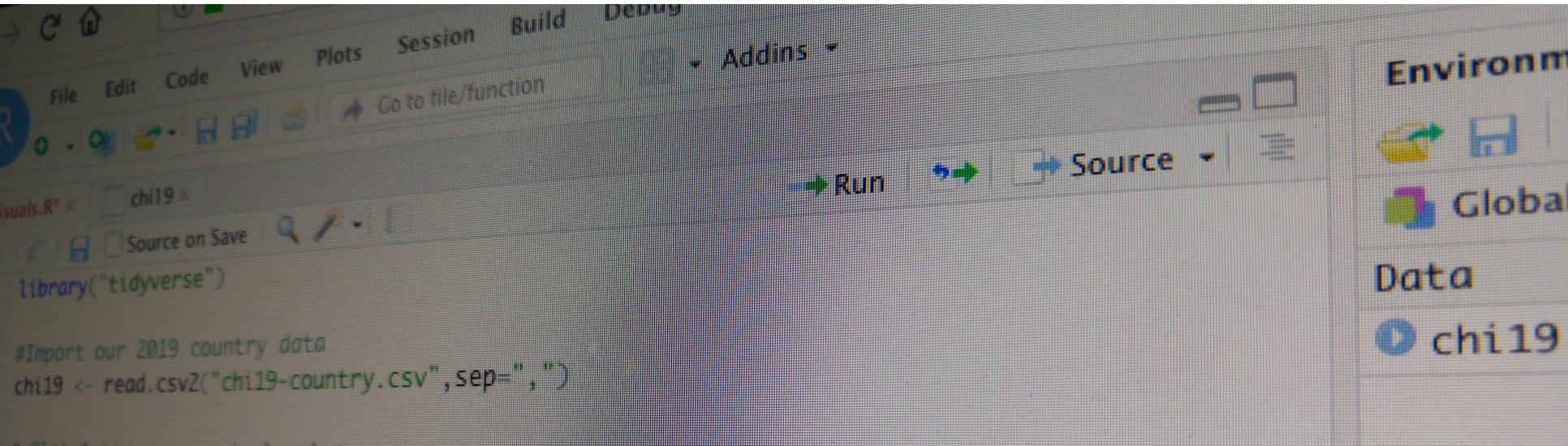
We're going to be using RStudio in the browser.

**<https://rstudio.sjg.uk>**

Log in with the details that I have given you.



Let's have a look at an example together



# Our dataset

We're going to be visualizing some CHI data.

They're data about how publications by country in 2018 and 2019.

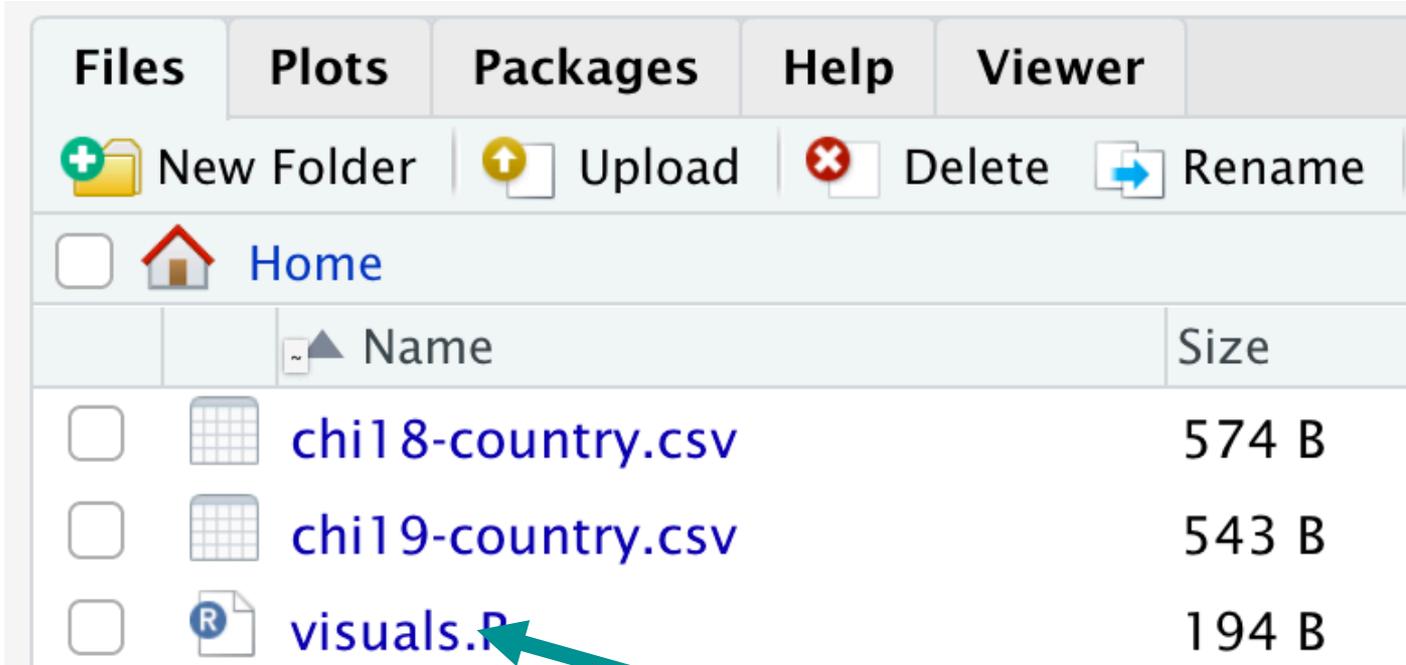
The 2019 data is based on Kashyap Todi's processing of data provided by the CHI Data Chair, Frederik Brudy.

The 2018 data are based on Kashyap Todi's data.

*Thanks to both!*



# Loading CHI 2019 data



The screenshot shows a file manager interface with a menu bar containing 'Files', 'Plots', 'Packages', 'Help', and 'Viewer'. Below the menu bar are icons for 'New Folder', 'Upload', 'Delete', and 'Rename'. The current directory is 'Home'. A table lists the files in the directory:

	Name	Size
<input type="checkbox"/>	 chi18-country.csv	574 B
<input type="checkbox"/>	 chi19-country.csv	543 B
<input type="checkbox"/>	 visuals.R	194 B

A teal arrow points from the bottom right towards the 'visuals.R' file name.

**Click this!**

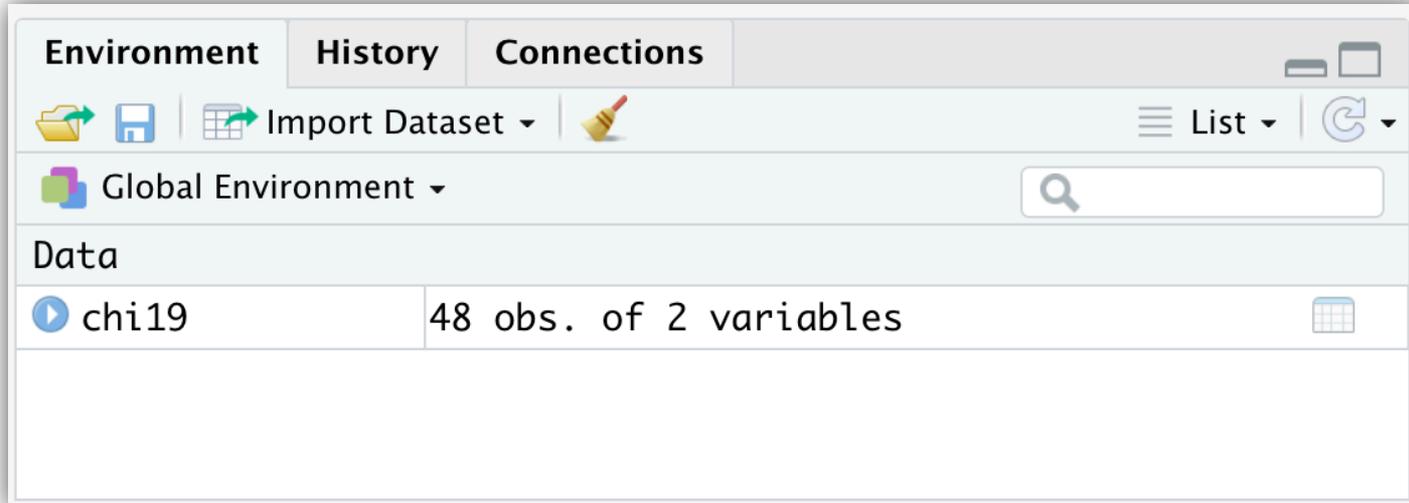


# CSV Import code

```
library("tidyverse")
```

```
#Import our 2019 country data
```

```
chi19 <- read.csv2("chi19-country.csv", sep=",")
```



The screenshot shows the RStudio Environment pane. At the top, there are tabs for 'Environment', 'History', and 'Connections'. Below the tabs is a toolbar with icons for file operations and a search bar. The main area displays the 'Global Environment' with a search bar. Under the 'Data' section, there is a table with one row: 'chi19' with '48 obs. of 2 variables'.

Data	
chi19	48 obs. of 2 variables



# CHI 2019 Country Data

n	Country	Pubs
1	United States	359
2	United Kingdom	125
3	Canada	74
4	Germany	68
5	Australia	41
6	Japan	32
7	China	25
8	Republic of Korea	23
9	Denmark	22
10	Netherlands	21



# ggplot2

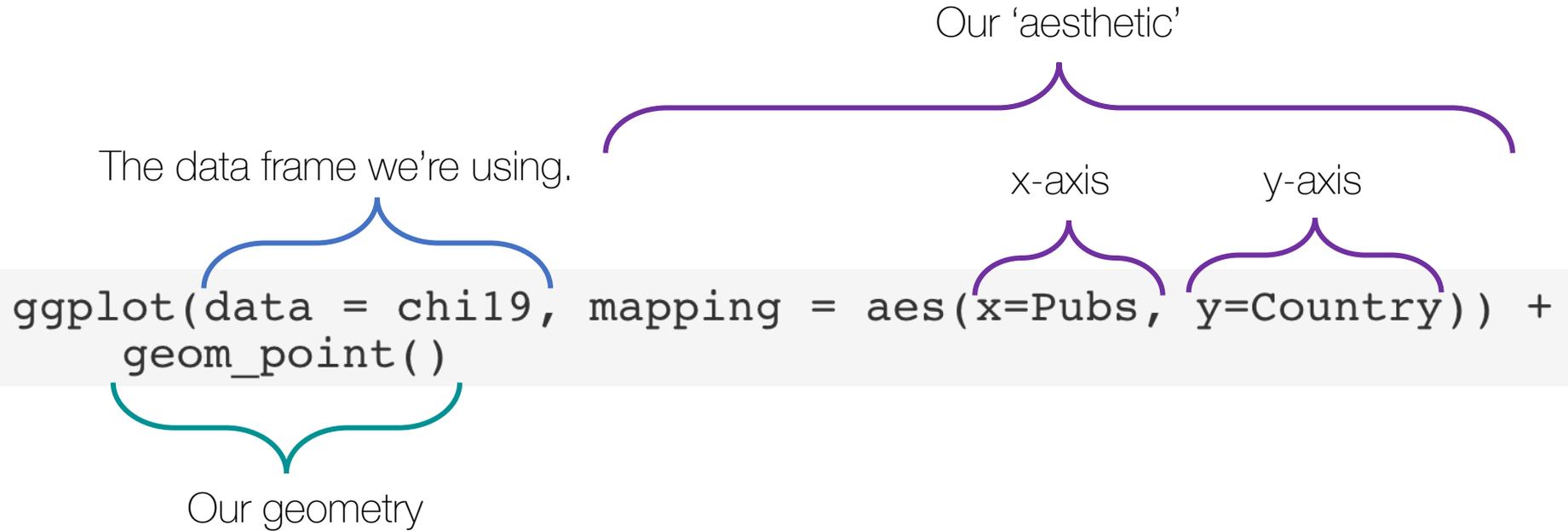
ggplot2 visualizations are created using the `ggplot()` function.

This function defines:

- The data we are using
- The aesthetics of the visualization we will create



# The ggplot2 function



# ggplot2 geometries



<code>geom_abline()</code>	<code>stat_sum()</code>	<code>geom_errorbar()</code>	<code>geom_segment()</code>
<code>geom_hline()</code>	<code>geom_density()</code>	<code>geom_linerange()</code>	<code>geom_curve()</code>
<code>geom_vline()</code>	<code>stat_density()</code>	<code>geom_pointrange()</code>	<code>geom_smooth()</code>
<code>geom_bar()</code>	<code>geom_density_2d()</code>	<code>geom_map()</code>	<code>stat_smooth()</code>
<code>geom_col()</code>	<code>stat_density_2d()</code>	<code>geom_path()</code>	<code>geom_spoke()</code>
<code>stat_count()</code>	<code>geom_dotplot()</code>	<code>geom_line()</code>	<code>geom_label()</code>
<code>geom_bin2d()</code>	<code>geom_errorbarh()</code>	<code>geom_step()</code>	<code>geom_text()</code>
<code>stat_bin_2d()</code>	<code>geom_hex()</code>	<code>geom_point()</code>	<code>geom_raster()</code>
<code>geom_blank()</code>	<code>stat_bin_hex()</code>	<code>geom_polygon()</code>	<code>geom_rect()</code> <code>geom_tile()</code>
<code>geom_boxplot()</code>	<code>geom_freqpoly()</code>	<code>geom_qq_line()</code>	<code>geom_violin()</code>
<code>stat_boxplot()</code>	<code>geom_histogram()</code>	<code>stat_qq_line()</code>	<code>stat_ydensity()</code>
<code>geom_contour()</code>	<code>stat_bin()</code>	<code>geom_qq()</code> <code>stat_qq()</code>	<code>stat_sf()</code> <code>geom_sf()</code>
<code>stat_contour()</code>	<code>geom_jitter()</code>	<code>geom_quantile()</code>	<code>geom_sf_label()</code>
<code>geom_count()</code>	<code>geom_crossbar()</code>	<code>stat_quantile()</code>	<code>geom_sf_text()</code>
<code>geom_ribbon()</code>	<code>geom_area()</code>	<code>geom_rug()</code>	<code>coord_sf()</code>



# The ggplot2 function

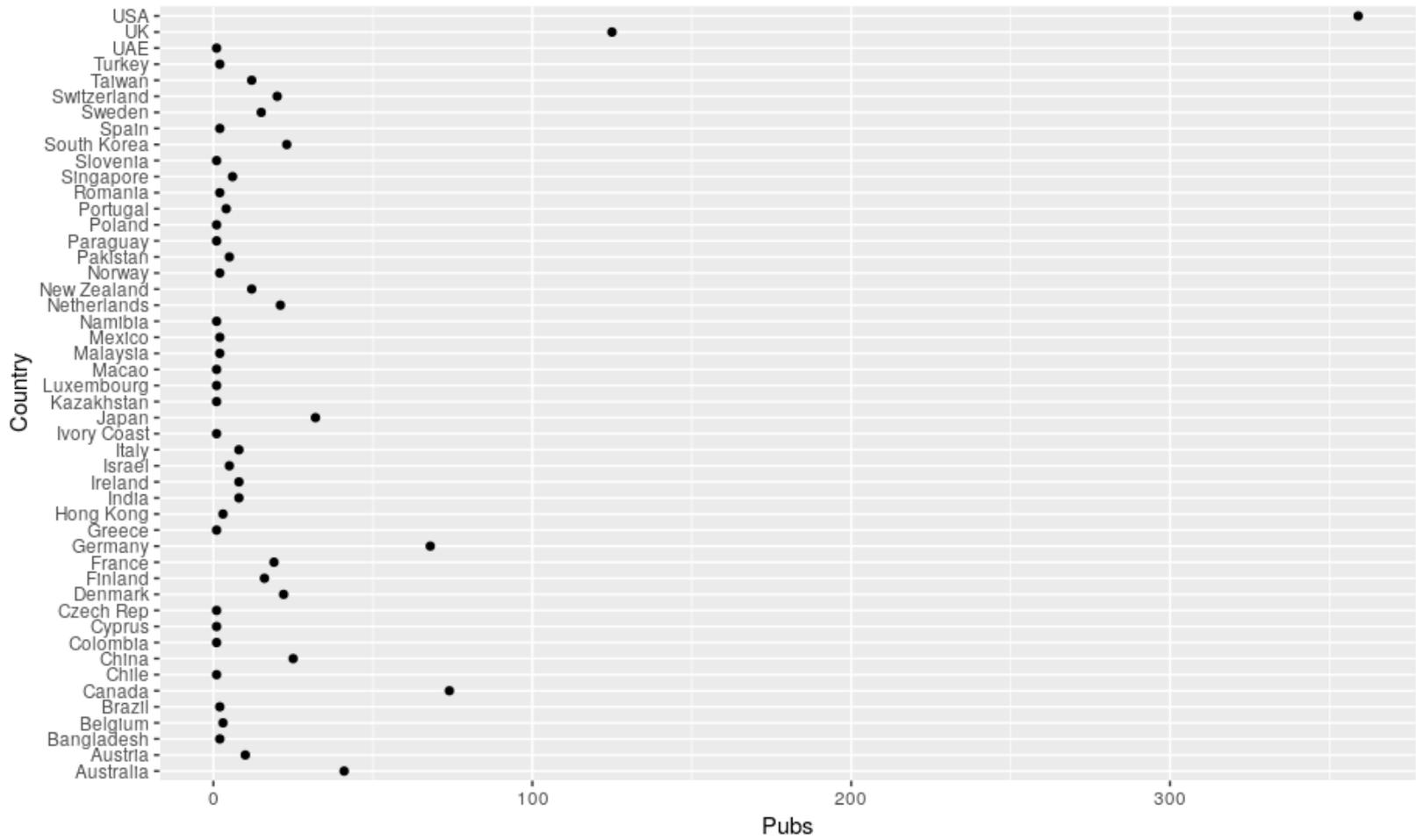
Now select the code and run it:

```
ggplot(data = chi19, mapping = aes(x=Pubs, y=Country)) +  
  geom_point()
```



**Click this!**





Let's try and build this out a bit

You should now have a basic plot. It's far from perfect though

Let's pick up the webpage at **5.3 Getting some order to proceedings**

You're going to:

- Get the data into a sensible order
- Have a play around with different geometries

Stop when you get to **6. Controlling the appearance of plots**



# Controlling the look of plots

There are three ways of controlling the look of plots:

– Through aesthetics

– Through geometries

– Through themes



We'll look at these first



# Appearance in geometries

```
ggplot(data = chi19, mapping = aes(x=reorder(Country, GDP), y=Year)) +  
  geom_point(colour="Orange", size=3) +  
  coord_flip()
```

Color of points

Size of points



# Appearance in geometries

On the website, work from

## **6.1 Appearance as part of geometries**

Until you reach the end of:

## **6.2 Controlling appearance with themes**



# Controlling appearance with themes

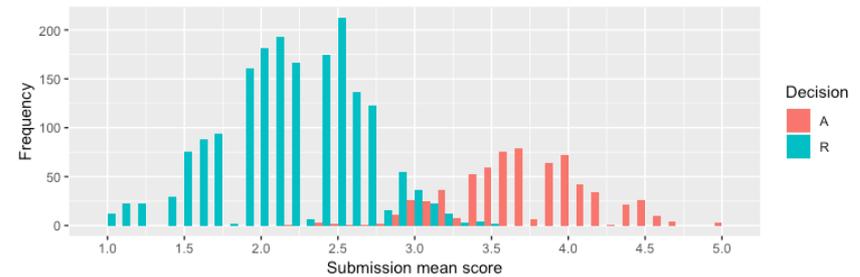
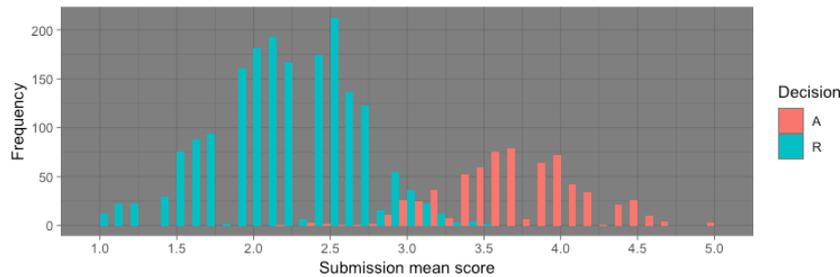
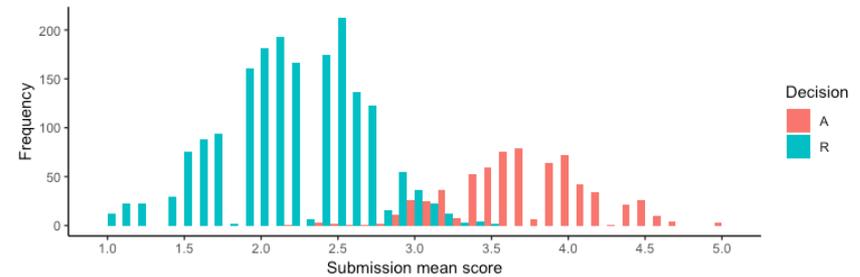
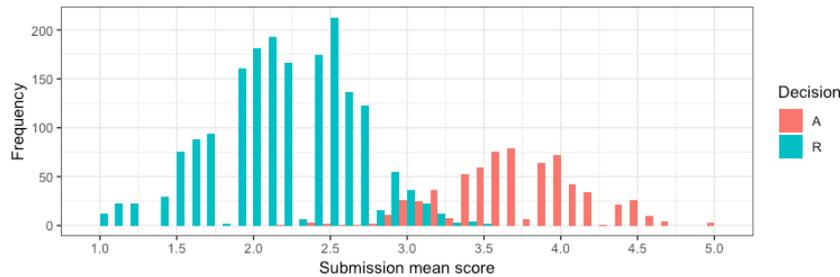
Themes control the **overall** appearance of your plot.

- Gridlines
- Background
- Axes
- Ticks
- Labels



# Controlling appearance with themes

There are built-in themes, like `theme_classic()`, `theme_dark()`, `theme_bw()` and `theme_gray()`.



# Controlling appearance with themes

Themes are made up of a number of components in a hierarchical form. (e.g., `legend.box.background` or `plot.caption`).

They are drawn using elements:

- Lines are `element_line()`
- Rectangles are `element_rect()`
- etc



# Controlling appearance with themes

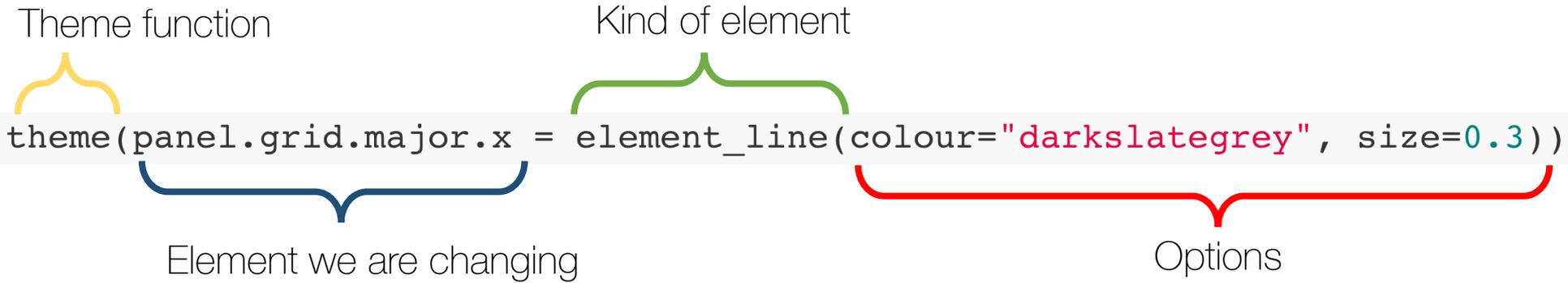
Theme function

Kind of element

```
theme(panel.grid.major.x = element_line(colour="darkslategrey", size=0.3))
```

Element we are changing

Options

The diagram shows the R code `theme(panel.grid.major.x = element_line(colour="darkslategrey", size=0.3))` with four annotations. A yellow bracket above the `theme` function is labeled "Theme function". A green bracket above `element_line` is labeled "Kind of element". A blue bracket below `panel.grid.major.x` is labeled "Element we are changing". A red bracket below `colour="darkslategrey", size=0.3` is labeled "Options".

# Controlling appearance with themes

On the website, work until you reach:

## **7. A plot drawing on two datasets**



# Plotting multiple variables

So far we have looked at a very simple plot.

- It has a categorical variable.
- And a continuous variable.

What if we want to look at many year's worth of data?

Enter the **CHI**mbined dataset with 2019 and 2018 data!



# CHI 2019 Country Data

<b>n</b>	<b>Country</b>	<b>Year</b>	<b>Pubs</b>
1	Australia	2018	23
2	Australia	2019	41
3	Austria	2018	8
4	Austria	2019	10
5	Belgium	2018	3
6	Belgium	2019	3
7	Brazil	2018	1
8	Brazil	2019	2
9	Canada	2018	78
10	Canada	2019	74



# Plotting multiple variables

There are three ways of controlling the look of plots:

- Through aesthetics 
- Through geometries
- Through themes



# Plotting multiple variables

aes function



```
aes(x=reorder(Country, Pubs), y=Pubs, colour=factor(Year))
```

New 'colour' aes addition



Stuff you recognize



# Plotting multiple variables

On the website, work from

## **7. A plot drawing on multiple variables**

Until you reach:

## **8. Mapping locations of CHI authors**

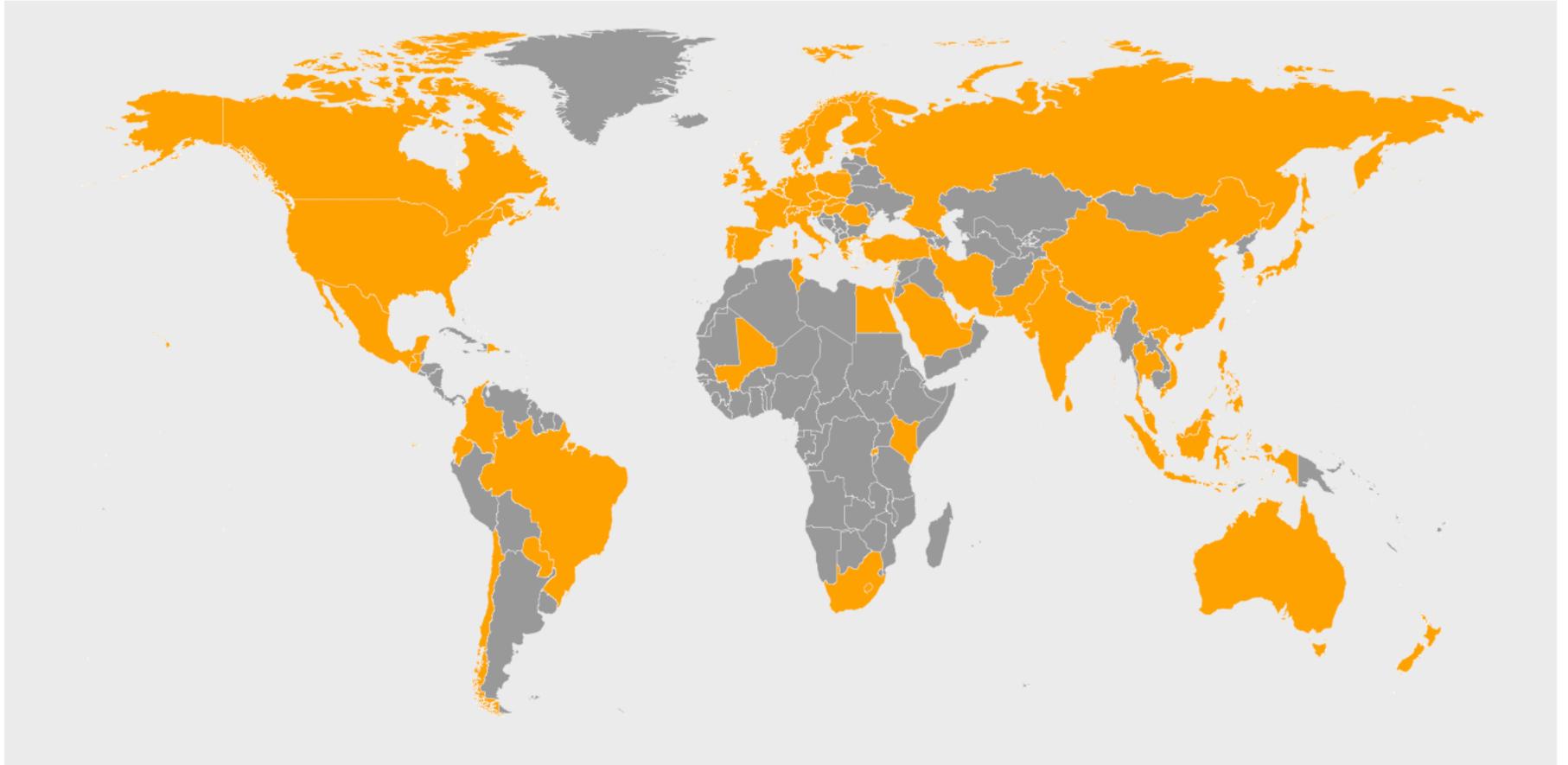


# Mapping locations

```
#Import mapping libraries  
library(maps)
```



# Mapping locations



# Mapping locations

```
world <- map_data("world")  
world$fac <- world$region %in% chi19$Country
```



# Mapping locations

We're going to use the `geom_poly()*` geometry for plotting our graph.

Map of the world      Longitude and latitude

```
ggplot(data = world, aes(x=long, y = lat, group = group)) +  
  geom_polygon(colour="white")
```

Plot the polygons      Colour borders in white      Keep countries together

\* There is a `geom_map()`, which I discuss on the website



# Mapping locations

We're going to use the `geom_poly()*` geometry for plotting our graph.

Map of the world      Longitude and latitude

```
ggplot(data = world, aes(x=long, y = lat, group = group)) +  
  geom_polygon(colour="white")
```

Plot the polygons      Colour borders in white      Keep countries together

\* There is a `geom_map()`, which I discuss on the website



# Plotting multiple variables

On the website, work from

## **8. Mapping locations of CHI authors**

Until you reach:

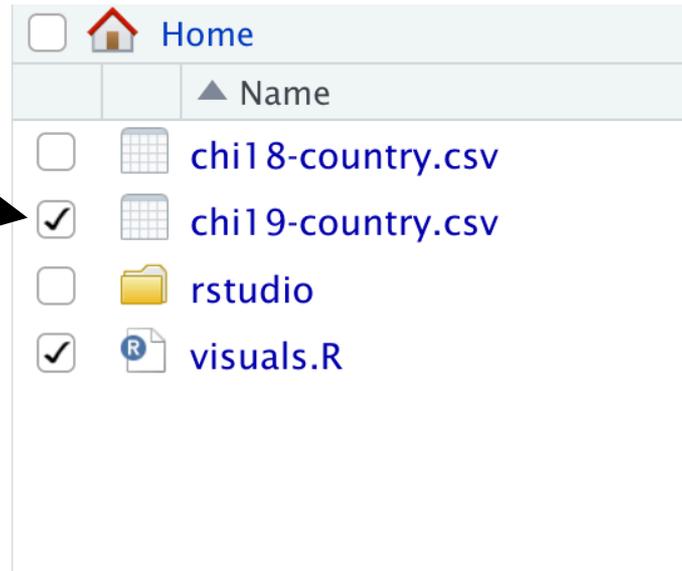
## **8. We're done... for now!**



# Getting your code

Download your files before you lose them!

**Select the  
files you want  
to keep**



# Getting your code

Download your files before you lose them!

**Then click  
More →  
Export...**

**To get a zip  
file.**

