

Human Geography through Merseyside - Quantitative Block: Seeing the world through numbers

Zi Ye and Ron Mahabir

2026-02-19

Table of contents

| | |
|---|-----------|
| Welcome | 4 |
| Contact | 4 |
| Overview | 5 |
| Aim and Learning Objectives | 5 |
| Module Structure | 5 |
| Software and Data | 6 |
| Assessment | 8 |
| 1 Lab: Getting Started in RStudio - Knowing Merseyside | 9 |
| 1.1 Overview | 9 |
| 1.2 Getting set up with RStudio | 9 |
| 1.2.1 Install R and RStudio (if necessary) | 9 |
| 1.2.2 File management | 10 |
| 1.2.3 Open RStudio | 10 |
| 1.2.4 Ways of working | 11 |
| 1.2.5 Your first R code | 12 |
| 1.3 Knowing Merseyside | 25 |
| 1.3.1 Merseyside districts | 25 |
| 1.3.2 Merseyside neighbourhoods | 41 |
| 1.4 Summary | 51 |
| 1.4.1 References | 51 |
| 1.5 Formative Tasks | 51 |
| 2 Lab: Exploratory Data Analysis - UK Election | 55 |
| 2.1 Overview | 55 |
| 2.2 Clear the decks | 55 |
| 2.3 Open libraries | 56 |
| 2.4 Parliamentary Constituency Data | 56 |
| 2.4.1 Load the dataset | 56 |
| 2.4.2 Familiar with the dataset and variable types | 57 |
| 2.4.3 Exploratory Data Analysis (EDA) | 58 |
| 2.5 Make your own map for the election result | 76 |
| 2.5.1 Read in Parliamentary Constituency Boundaries | 76 |
| 2.5.2 Inspect the spatial dataset | 77 |

| | | |
|----------|--|------------|
| 2.5.3 | Link boundaries to pc_data | 78 |
| 2.5.4 | Map the election result | 79 |
| 2.6 | Formative tasks | 80 |
| 3 | Lab: Introductory Statistics - Happiness around the world | 87 |
| 3.1 | Overview | 87 |
| 3.2 | Prepare your working environment | 87 |
| 3.3 | Load libraries and familiar with survey data | 88 |
| 3.4 | Sampling of numerical variables | 89 |
| 3.4.1 | Sample mean and Standard Error | 91 |
| 3.4.2 | Confidence Intervals | 94 |
| 3.5 | Sampling of categorical variable | 95 |
| 3.5.1 | Sample proportion and Standard Error of the proportion | 97 |
| 3.5.2 | Confidence Intervals for a proportion (95%) | 100 |
| 3.6 | Compare to World Value Survey | 101 |
| 3.6.1 | Happiness in different survey sample | 101 |
| 3.6.2 | Between regions and countries | 104 |
| 3.6.3 | Map happiness around the world | 108 |
| 3.7 | Formative Tasks | 111 |
| 4 | Lab: Correlation, data reliability and the issue of scale - Health | 118 |
| 4.1 | Overview | 118 |
| 4.2 | Prepare your working environment | 118 |
| 4.3 | Load libraries | 119 |
| 4.4 | The Geography of Health (area-level) | 119 |
| 4.4.1 | Viewing the district level boundaries | 119 |
| 4.4.2 | The spatial distribution Poor Health | 122 |
| 4.4.3 | Visualise associations with a scatter plot | 126 |
| 4.5 | Measuring the correlation of <i>continuous</i> variables | 126 |
| 4.5.1 | Check if both variables are symmetrical | 127 |
| 4.5.2 | Calculate appropriate measure of correlation | 129 |
| 4.5.3 | Visualising the association between categorical variables | 130 |
| 4.6 | Measuring the correlation of categorical variables | 132 |
| 4.7 | Area-level data analysis traps - MAUP | 134 |
| 4.8 | Formative tasks | 136 |

Welcome

This is the website for “Human Geography through Merseyside - Quantitative Block: Seeing the world through numbers” (module **ENVS162**) at the University of Liverpool. This block of the module is designed and delivered by Dr. Zi Ye and Dr. Ron Mahabir from the Geographic Data Science Lab at the University of Liverpool. The module seeks to provide hands-on experience and training in introductory statistics for human geographers.

The website is **free to use** and is licensed under the [Attribution-NonCommercial-NoDerivatives 4.0 International](#). A compilation of this web course is hosted as a GitHub repository that you can access:

- As an [html website](#).
- As a [GitHub repository](#).

Contact

Zi Ye - zi.ye [at] liverpool.ac.uk Lecturer in Geographic Information Science Office 107, Roxby Building, University of Liverpool - 74 Bedford St S, Liverpool, L69 7ZT, United Kingdom.

Ron Mahabir - Ron.Mahabir [at] liverpool.ac.uk Lecturer in Geographic Data Science Office 4xx, Roxby Building, University of Liverpool - 74 Bedford St S, Liverpool, L69 7ZT, United Kingdom.

Overview

Aim and Learning Objectives

This sub-module aims to provide training and skills on a set of basic quantitative skills for data collection, analysis, and interpretation and to enable you to link conceptual ideas with real world examples. **This block serves as the foundation for Year 2 BA field class and, optionally, for Year 3 dissertation.**

Background

Data and research are key pillars of the global economy and society today. We need rigorous approaches to collecting and analysing both the statistics that can tell us ‘how much’ and if there are observable relationships between phenomena; and the information gives us a nuanced understanding of cultural contexts and human dynamics. Quantitative skills enable us to explore and measure socio-economic activities and processes at large scales, while qualitative skills enable understanding of social, cultural, and political contexts and diverse lived experiences. Rather than being in opposition, qualitative and quantitative research can complement one another in the investigation of today’s pressing research questions.

To these ends, this block will help you develop your quantitative skills, as critical tools. This course will help you understand what quantitative statistical researchers use and develop a set of research techniques that can be used in your field classes and dissertations.

Learning objectives:

- Understand how to explore a dataset, containing a number of observations described by a set of variables.
- Demonstrate an understanding in the application and interpretation of commonly used quantitative research methods.
- Ability to work with quantitative data to understand real-world social phenomenon and patterns.

Module Structure

Staff: Dr Zi Ye and Dr Ron Mahabir

Where and When

Week 1 - 5 Lecture: Tuesday (12am – 1pm) @ Mathematical Sciences, Proudman Lecture Theatre

Week 1 - 6 Practical PC session: Friday (9 – 11 am) @ Central Teaching Lab: PC Teaching Centre

Lectures will introduce and explain the fundamentals of quantitative methods, with the opportunity to apply the method introduced in the labs later in the week.

The computer practical sessions, will give you the chance to use and apply quantitative methods to real-world data. These are primarily self-directed sessions, but with support on hand if you get stuck. Support and training in R will be provided through these sessions. Weekly sessions will be driven by empirical research questions.

| Week | Topic | Format | Staff |
|------|--|------------------------------------|-------|
| 1 | Introduction Getting Started in RStudio: Knowing Merseyside | Lecture Computer Lab Practical | ZY/RM |
| 2 | Exploratory Data Analysis: UK Election | Lecture and Computer Lab Practical | ZY |
| 3 | Sampling and data manipulation: Happiness around the world | Lecture and Computer Lab Practical | ZY |
| 4 | Correlation, data reliability and the issue of scale: Health | Lecture and Computer Lab Practical | RM |
| 5 | How robust are my findings | Lecture and Computer Lab Practical | RM |
| 6 | Online Assessment | Computer Lab | RM/ZY |

Software and Data

For quantitative training sessions, ensure you have installed and/or have access to **RStudio**. To run the analysis and reproduce the code in R, you need the following software installed on your machine:

- R-4.2.2 (or later)
- RStudio 2022.12.0-353 (or later)

To install and update:

- R, download the appropriate version from [The Comprehensive R Archive Network \(CRAN\)](#).
- RStudio, download the appropriate version from [here](#).

This software is already installed on University Machines. But you will need it to run the analysis on your personal devices.

Data

Example datasets could be accessed through [ENVS 162 Canvas module](#) every week.

Assessment

Week 6 Computer-based ‘open book’ multiple-choice exam

- The online assessment will be released at **4pm on Thursday 5th March** and should be **completed by 4pm on Friday 6th March**.
- **Also available 06/03/2025 9:00 – 11:00 CTL PC Teaching Centre (1st Floor CTL)**
- Should take less **90 minutes**; c. 20 questions; 24 hours to complete
- Questions and answers randomised for each student (anti-cheating measure)
- Some questions of factual recall, more requiring data analysis to find answers

Preparation for assessment

- Weekly lecture & weekly computer practical ‘clinic sessions’
- Weekly holding hands formative tasks at the last 20 mins of the practical session
- Week 5 mock online test

1 Lab: Getting Started in RStudio - Knowing Merseyside

1.1 Overview

This practical intend to prepare students who have limited experiences with R and RStudio. The content are adapted based on

- Brunsdon, Chris, and Lex Comber. 2018. *An Introduction to r for Spatial Analysis and Mapping (2e)*. Sage.
- Comber, Lex, and Chris Brunsdon. 2021. *Geographical Data Science and Spatial Data Analysis: An Introduction in r*. Sage.

1.2 Getting set up with RStudio

1.2.1 Install R and RStudio (if necessary)

R is a free, open-source programming language used for statistical analysis, data visualization, and data science

RStudio is a free front-end to R, designed to make using R easier

All of the PCs in the University PC Teaching Centre used for this class come with R and RStudio pre-installed, as do the PCs in many other University PC Teaching Centres.

However, you may wish to install R and RStudio on your own computer, or on a University PC that lacks them.

University computers: Use the *Install University Applications* app on the computer to install the latest version of RStudio (this will also install the latest version of R)

Your own computer: R and RStudio can be downloaded from the CRAN website and installed your own computer - see below for details. **A key point is that you should install R before you install RStudio.**

The simplest way to get R installed on your computer is to go the download pages on the R website - a quick search for 'download R' should take you there, but if not you could try:

- Windows: <https://cran.r-project.org/bin/windows/base/>
- Mac: <https://cran.r-project.org/bin/macosx/>
- Linux: <http://cran.r-project.org/bin/linux/>

The Windows and Mac version come with installer packages and are easy to install whilst the Linux binaries require use of a command terminal.

RStudio can be downloaded from <https://www.rstudio.com/products/rstudio/download/> and the free version of RStudio Desktop is more than sufficient for this module and all the other things you will to do at degree level.

If you experience any problems installing R or RStudio on your own computer, bring it to one of the class lab sessions where we will be able to provide advice.

1.2.2 File management

Before you start installing software or downloading data, create a folder on your M-Drive (if working on a University networked machine) or locally if working on your own device – name this ‘ENVS162’ and within this create a sub-folder for each practical session. For this session, create a sub-folder called **Week1** in your **ENVS162** folder on your M-Drive. Take care to ensure you do not delete any work you do complete in the practical sessions. It is imperative that you practice good file management!

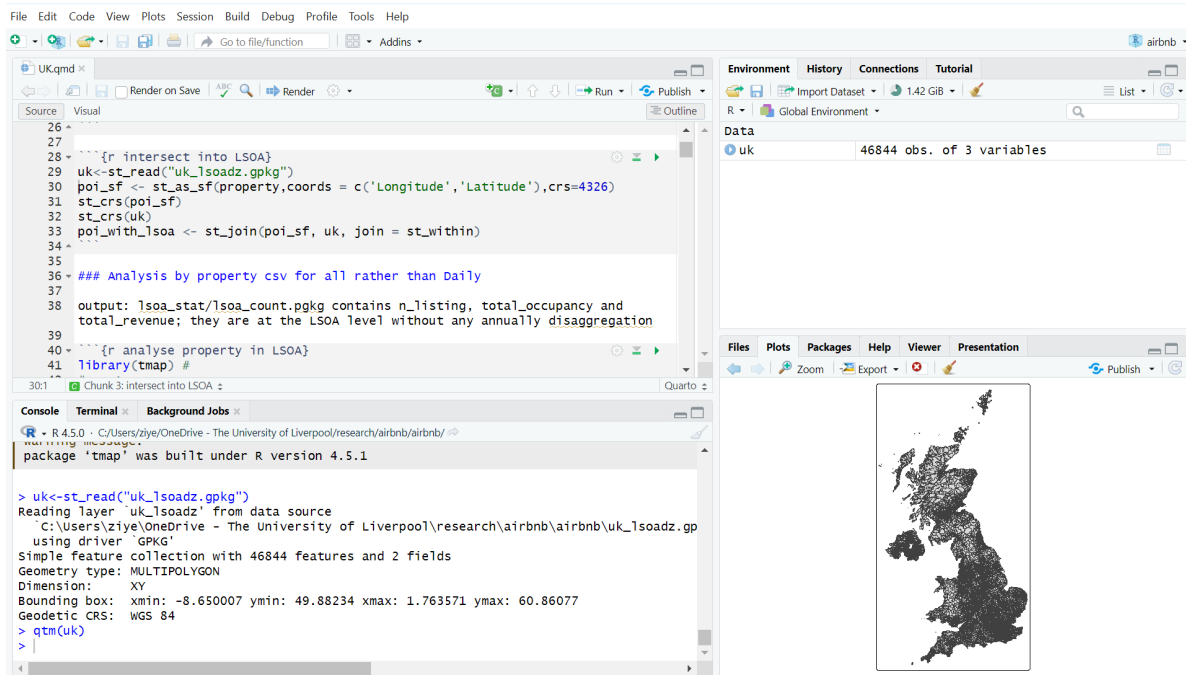
1.2.3 Open RStudio

RStudio provides an interface to the different things that R can do via the 4 panes: the Console where code is entered (bottom left), a Source pane with R scripts (top left), the variables in the working Environment (top right), Files, Plots, Help etc (bottom right) - see the RStudio environment in Figure below.

In the figure above of the RStudio interface, a new script has been opened, a line of code had been written and then run in the console. The code assigns a value of 100 to **x**. The file has been saved into the current working environment. You are expected to define a similar set up for each practical as you work through the code. Note that **in the script**, anything that follows a **#** is a comment and ignored by R.

Users can set up their personal preferences for how they like their RStudio interface. Similar to straight R, there are very few pull-down menus in R, and therefore you will type lines of code into your script and run these in what is termed a *command line interface* (the console). Like all command line interfaces, the learning curve is steep but the interaction with the software is more detailed which allows greater flexibility and precision in the specification of commands.

Beyond this there are further choices to be made. Commands can be entered in two forms: directly into the *R console* window or as a series of commands into a script window. We strongly advise that all code should be **written in a script** - (a *.R* file) and then **run from the script**. - To run code in a script, place the cursor on the line of code and then run by pressing the 'Run' icon at the top left of the script pane, or by pressing **Ctrl Enter** (PC) (or **Cmd Enter** on a Mac).



1.2.4 Ways of working

The first set of consideration relate to *how* you should work in R/RStudio. The key things to remember are:

- R is a learning curve if you have never done anything like this before. It can be scary. It can be intimidating. But once you have a bit of familiarity with how things work, it is incredibly powerful.
- You will be working from practical worksheets which will have all the code you need. Your job is to try to **understand** what the code is doing and **not** to remember the code. Comments in your code really help.
- To help you do this, the very strong suggestion is use the R scripts that are provided, and that you add your own comments to help you understand what is going on when you return to them. Comments are prefaced by a hash (#) that is ignored by R. Then

you can save your code (with comments), run it and return to it later and modify at your leisure.

The module places a strong emphasis placed on learning by doing, which means that you are encouraged to unpick the code that you are given, adapt it and play with it. It is not about remembering or being able to recall each function used but about understanding what is being done. If you can remember what you did previously (i.e. the operations you undertook) and understand what you did, you will be able to return to your code the next time you want to do something similar. To help you with this you should:

1. Always run your code from an R script... **always!** These are provided for each practical;
2. Annotate your scripts with comments. These are prefixed by a hash (#) in the code;
3. Save your R script to your folder.

In Summary:

- You should always use a script (a text file containing code) for your code which can be saved and then re-run at a later date.
- You can write your own code into a script, copy and paste code into it or use an existing script (for example as provided for each of the R/RStudio practicals in this module).
- To open a new R script go to File > New File > R Script to open a new R file, and save it with a sensible name.
- To load an existing script file go to File > Open File and then navigate to your file. Or, if you have recently opened the file, go to File > Recent Files >.
- It is good practice to set the working directory at the beginning of your R session. This can be done via the menu in RStudio Session > Set Working Directory > This points the R session to the folder you choose and will ensure that any files you wish to read, write or save are placed in this directory.
- To run code in a script, place the cursor on the line of code and then run by pressing the 'Run' icon at the top left of the script pane, or by pressing Ctrl Enter (PC) or Cmd Enter (Mac).

1.2.5 Your first R code

In this section you will undertake a few generic operations. You will:

- undertake **assignment**: the allocation of values to an R object.
- use assignment to create a **vector** of elements and a **matrix** of elements.
- undertake **operations** on R objects.

- apply some **functions** to R objects (functions nearly always return a value).
- access some of R in-built data to examine a data table (or **data.frame** which is like an Excel spreadsheet).
- do some basic **plotting**, including scatter plots and histograms.
- create data summaries.

On the way you will also be introduced to **indexing**.

First, you should **create a new R script** (see above) and save it as **week1.R** in the working directory you are using for this practical. This should be the **Week1** sub-directory you created in the **ENVS162** folder. Note that you should create a separate folder for each week's practical.

1.2.5.1 Assignment

The command line prompt in the Console window, the **>**, is an invitation to start typing in your commands.

Write the following into your script: **3+5** and run it. Recall that code is run done by either by pressing the Run icon at the top left of the script pane, or by pressing **Ctrl Enter** (PC) or **Cmd Enter** (Mac).

```
3+5
```

```
[1] 8
```

Here the result is 8. The **[1]** that precedes the output it formally indicates, *first requested element will follow*. In this case there is just one element. The **>** indicates that R is ready for another command.

Now type the following in to your script and run it:

```
y <- 3+5
y
```

```
[1] 8
```

Here the value of the 3+5 has been *assigned* to y. The syntax `y <- 3+5` can be read as y *gets* 3+5. When y is invoked its value is returned (8).

For the purposes of this module, in R the equals sign (=) is the same as <-, a left diamond bracket < followed by a minus sign -. This too is interpreted by R as *is assigned to* or *gets* when the code is read **right to left**.

Now copy and paste the following into your R script and run both lines:

```
x <- matrix(c(1,2,3,4,5,6,7,8), nrow = 4)
y = matrix(1:8, nrow = 4, byrow = T)
```

You should see the x appear with the y in the Environment pane. y has now been overwritten with a new assignment. If you click on the icon next to them, you will get a ‘spreadsheet’ view of the data you have created.

Of course you can also enter the following in the console and see what is returned:

x

| | [,1] | [,2] |
|------|------|------|
| [1,] | 1 | 5 |
| [2,] | 2 | 6 |
| [3,] | 3 | 7 |
| [4,] | 4 | 8 |

y

| | [,1] | [,2] |
|------|------|------|
| [1,] | 1 | 2 |
| [2,] | 3 | 4 |
| [3,] | 5 | 6 |
| [4,] | 7 | 8 |

Note In the code snippets above you have used **parentheses** - round brackets. Different kinds of brackets are used in different ways in R. Parentheses are used with **functions**, and contain the **arguments** that are passed to the function, separated by commas (,).

In this case the functions are `c()` and `matrix()`. The function `c()` combines or concatenates elements into a vector, and `matrix()` creates a matrix of elements in a tabular format.

In the line of code `x = matrix(c(1,2,3,4,5,6,7,8), nrow = 4)`, the arguments passed to the `matrix()` function are the vector of values `c(1,2,3,4,5,6,7,8)` and `nrow = 4`. Other kinds of brackets are used in different ways as you will see later.

One final thing to note is that the matrix `x` has the numbers 1 to 8, but this is specified by `1:8`. Try entering `3:65`, `19:11`, and `1.5:5` to see how the colon (`:`) works in this context.

1.2.5.2 Operations

Now you can undertake *operations* on R objects and apply *functions* to them. Write the following code into your script and then run it:

```
# x is a matrix
x
```

```
      [,1] [,2]
[1,]     1     5
[2,]     2     6
[3,]     3     7
[4,]     4     8
```

```
# multiplication
x*2
```

```
      [,1] [,2]
[1,]     2    10
[2,]     4    12
[3,]     6    14
[4,]     8    16
```

```
# sum of x
sum(x)
```

```
[1] 36
```

```
# mean of x
mean(x)
```

```
[1] 4.5
```

Operations can be undertaken directly using mathematical notation like `*` for multiplication or using functions like `max` to find the maximum value in an R object.

1.2.5.3 Functions

Functions are always followed by parenthesis (round brackets) (). These are different from square and curly brackets [] and { }. Functions always return something, a result if you like, and have the generic form:

```
# don't run this or write this into your script!
result = function(value or R object, other parameters)
```

Do not run or enter this code in your script - it is an example!

1.2.5.4 Data Tables

Here we will load a data table in `data.frame` (like a spreadsheet) in R/RStudio. R has number of in-built datasets that we can use the code below loads one of these:

```
data(mtcars)
class(mtcars)
```

```
[1] "data.frame"
```

Have a look at what is loaded by listing the objects in the current R session

```
ls()
```

```
[1] "mtcars" "x"      "y"
```

You should see the `mtcars` object. You can examine this data in a number of ways

```
# the structure of mtcars
str(mtcars)
```

```
'data.frame':  32 obs. of  11 variables:
 $ mpg : num  21 21 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 ...
 $ cyl : num   6  6  4  6  8  6  8  4  4  6 ...
 $ disp: num  160 160 108 258 360 ...
 $ hp  : num  110 110  93 110 175 105 245  62  95 123 ...
 $ drat: num   3.9 3.9 3.85 3.08 3.15 2.76 3.21 3.69 3.92 3.92 ...
 $ wt  : num   2.62 2.88 2.32 3.21 3.44 ...
```



```
$ qsec: num 16.5 17 18.6 19.4 17 ...
$ vs : num 0 0 1 1 0 1 0 1 1 1 ...
$ am : num 1 1 1 0 0 0 0 0 0 0 ...
$ gear: num 4 4 4 3 3 3 3 4 4 4 ...
$ carb: num 4 4 1 1 2 1 4 2 2 4 ...
```

```
# the first six rows (or head) of mtcars
head(mtcars)
```

| | mpg | cyl | disp | hp | drat | wt | qsec | vs | am | gear | carb |
|-------------------|------|-----|------|-----|------|-------|-------|----|----|------|------|
| Mazda RX4 | 21.0 | 6 | 160 | 110 | 3.90 | 2.620 | 16.46 | 0 | 1 | 4 | 4 |
| Mazda RX4 Wag | 21.0 | 6 | 160 | 110 | 3.90 | 2.875 | 17.02 | 0 | 1 | 4 | 4 |
| Datsun 710 | 22.8 | 4 | 108 | 93 | 3.85 | 2.320 | 18.61 | 1 | 1 | 4 | 1 |
| Hornet 4 Drive | 21.4 | 6 | 258 | 110 | 3.08 | 3.215 | 19.44 | 1 | 0 | 3 | 1 |
| Hornet Sportabout | 18.7 | 8 | 360 | 175 | 3.15 | 3.440 | 17.02 | 0 | 0 | 3 | 2 |
| Valiant | 18.1 | 6 | 225 | 105 | 2.76 | 3.460 | 20.22 | 1 | 0 | 3 | 1 |

The `mtcars` object is a `data.frame`, a kind of data table, and it has a number of attributes which are all numeric. The code below prints it all out to the console:

```
mtcars
```

Data frames are ‘flat’ in that they typically have a rectangular layout like a spreadsheet, with rows typically relating to observations (individuals, areas, people, houses, etc) and columns relating to their properties or attributes (height, age, etc). The columns in data frames can be of different types: vectors of numbers, factors (classes) or text strings. In matrices all of the columns have to be of the same type. Data frames are central to what we will do in R.

1.2.5.5 Plotting the data: ‘Hello World!’

The code below creates a plot of 2 variables counts in the data: `mpg` and `disp`.

```
plot(disp ~ mpg, data = mtcars, pch=16)
```



The option `pch=16` sets the plotting character to a solid black dot. More plot characters are available - examine the help for `points()` to see these (For any command, if you are the first time use it, you can always ask R to explain to you by using `?` as help)

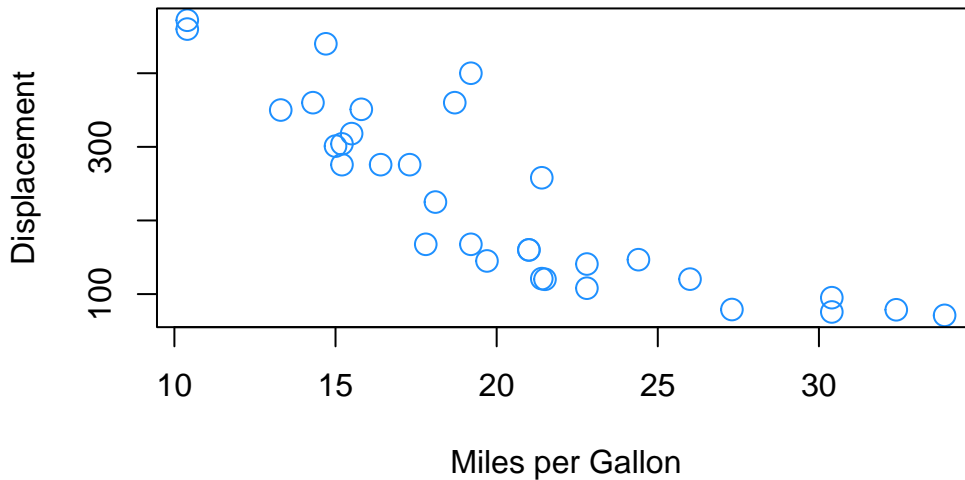
```
?points
```

This plot can be improved greatly. We can specify more informative axis labels, change size of the text and of the plotting symbol, and so on.

We can also specify the same plot by passing named variables to the `plot` function directly as well as other parameters, as in the figure. Notice how the syntax for this is different to the `plot` function above, and the different **parameters** that are passed to the `plot()` function:

```
plot(x = mtcars$mpg, y = mtcars$disp,   pch = 1, col = "dodgerblue",
     cex = 1.5, xlab = "Miles per Gallon", ylab = "Displacement",
     main = "Hello World!")
```

Hello World!



Notice how the dollar sign (\$) is used to access variables in the `mtcars` data table compared to the first plot command, which specified `data = mtcars`.

1.2.5.6 Data summaries and indexing

We may for example require information on variables in `mtcars`. The `summary` function is very useful:

```
summary(mtcars)
```

| mpg | cyl | disp | hp |
|---------------|---------------|---------------|---------------|
| Min. :10.40 | Min. :4.000 | Min. : 71.1 | Min. : 52.0 |
| 1st Qu.:15.43 | 1st Qu.:4.000 | 1st Qu.:120.8 | 1st Qu.: 96.5 |
| Median :19.20 | Median :6.000 | Median :196.3 | Median :123.0 |
| Mean :20.09 | Mean :6.188 | Mean :230.7 | Mean :146.7 |
| 3rd Qu.:22.80 | 3rd Qu.:8.000 | 3rd Qu.:326.0 | 3rd Qu.:180.0 |
| Max. :33.90 | Max. :8.000 | Max. :472.0 | Max. :335.0 |

| drat | wt | qsec | vs |
|---------------|---------------|---------------|----------------|
| Min. :2.760 | Min. :1.513 | Min. :14.50 | Min. :0.0000 |
| 1st Qu.:3.080 | 1st Qu.:2.581 | 1st Qu.:16.89 | 1st Qu.:0.0000 |
| Median :3.695 | Median :3.325 | Median :17.71 | Median :0.0000 |
| Mean :3.597 | Mean :3.217 | Mean :17.85 | Mean :0.4375 |

| | | | |
|----------------|---------------|---------------|----------------|
| 3rd Qu.:3.920 | 3rd Qu.:3.610 | 3rd Qu.:18.90 | 3rd Qu.:1.0000 |
| Max. :4.930 | Max. :5.424 | Max. :22.90 | Max. :1.0000 |
| am | gear | carb | |
| Min. :0.0000 | Min. :3.000 | Min. :1.000 | |
| 1st Qu.:0.0000 | 1st Qu.:3.000 | 1st Qu.:2.000 | |
| Median :0.0000 | Median :4.000 | Median :2.000 | |
| Mean :0.4062 | Mean :3.688 | Mean :2.812 | |
| 3rd Qu.:1.0000 | 3rd Qu.:4.000 | 3rd Qu.:4.000 | |
| Max. :1.0000 | Max. :5.000 | Max. :8.000 | |

This shows different summaries of the individual attributes in `mtcars`.

The main R graphics function is `plot()`. In addition to `plot()` there are functions for adding points and lines to existing graphs, for placing text at specified positions, for specifying tick marks and tick labels, for labelling axes, and so on.

There are various other alternative helpful forms of graphical summary. A helpful graphical summary for the `mtcars` data frame is the scatterplot matrix.

```
# return the names of the mtcars variables
names(mtcars)
```

```
[1] "mpg" "cyl" "disp" "hp" "drat" "wt" "qsec" "vs" "am" "gear"
[11] "carb"
```

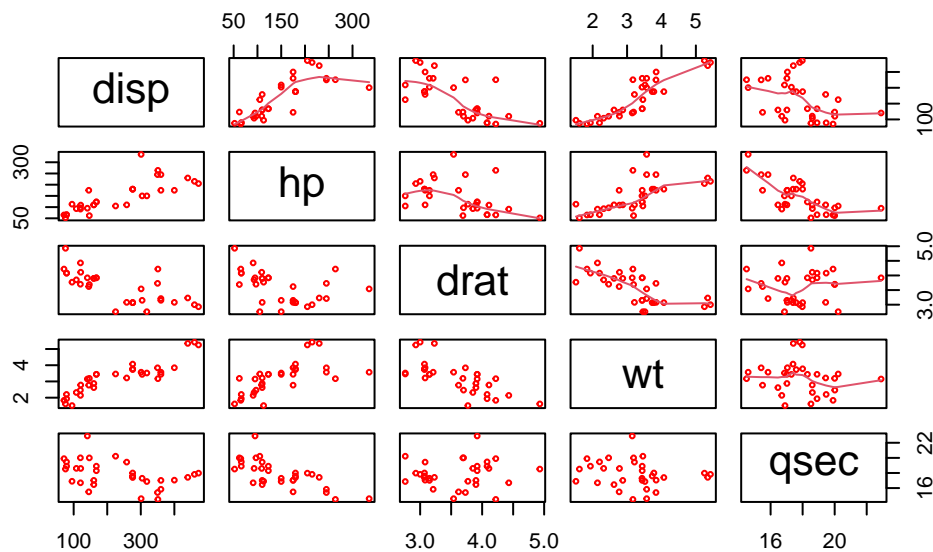
```
# return the 3rd to 7th names
names(mtcars)[c(3:7)]
```

```
[1] "disp" "hp" "drat" "wt" "qsec"
```

```
# check what this does
c(3:7)
```

```
[1] 3 4 5 6 7
```

```
# plot the 3rd to 7th variables in mtcars
plot(mtcars[, c(3:7)], cex = 0.5,
     col = "red", upper.panel=panel.smooth)
```



The results show the correlations between the variables in the `mtcars` data frame, and the trend of their relationship is included with the `upper.panel=panel.smooth` parameter passed to `plot`.

There are number of things to notice here (as well as the figure). In particular note the use of the vector `c(2:7)` to subset the columns of `mtcars`:

- In the second line, this is was used to subset the vector of column names created by `names(mtcars)`.
- In the third line, it was printed out. Notice how `3:7` printed out all the number between 3 and 7 - very useful.
- For the plot, the vector was passed to the second argument, after the comma, in the square brackets `[,]` to indicate which columns were to be plotted.

The referencing in this way (or *indexing*) is **very important**: the individual rows and columns of 2 dimensional data structures like data frames, matrices, tibbles etc can be accessed by passing references to them in the square brackets.

```
# 1st row
mtcars[1,]
```

```
      mpg  cyl  disp  hp  drat   wt   qsec vs  am  gear  carb
Mazda RX4   21    6  160 110  3.9 2.62 16.46  0   1    4     4
```

```
# 3rd column
mtcars[,3]
```

```
[1] 160.0 160.0 108.0 258.0 360.0 225.0 360.0 146.7 140.8 167.6 167.6 275.8
[13] 275.8 275.8 472.0 460.0 440.0  78.7  75.7  71.1 120.1 318.0 304.0 350.0
[25] 400.0  79.0 120.3  95.1 351.0 145.0 301.0 121.0
```

```
# a selection of rows
mtcars[c(3:5,8),]
```

| | mpg | cyl | disp | hp | drat | wt | qsec | vs | am | gear | carb |
|-------------------|------|-----|-------|-----|------|-------|-------|----|----|------|------|
| Datsun 710 | 22.8 | 4 | 108.0 | 93 | 3.85 | 2.320 | 18.61 | 1 | 1 | 4 | 1 |
| Hornet 4 Drive | 21.4 | 6 | 258.0 | 110 | 3.08 | 3.215 | 19.44 | 1 | 0 | 3 | 1 |
| Hornet Sportabout | 18.7 | 8 | 360.0 | 175 | 3.15 | 3.440 | 17.02 | 0 | 0 | 3 | 2 |
| Merc 240D | 24.4 | 4 | 146.7 | 62 | 3.69 | 3.190 | 20.00 | 1 | 0 | 4 | 2 |

Such indexing could of course have been assigned to a R object and used to do the subsetting:

```
x = c(3:7)
names(mtcars)[x]
```

```
[1] "disp" "hp"    "drat" "wt"    "qsec"
```

Thus indexing allows specific rows and columns to be extracted from the data as required.

Note You have encountered a second type of brackets, square brackets []. These are used to reference or **index** positions in a vector or a data table.

Consider the object `x` above. It contains a vector of values, 3,4,5,6,7. Entering `x[1]` would extract the first element of `x`, in this case 3. Similarly `x[4]` would return the 4th element and `x[c(1,4)]` would return the 1st and 4th elements of `x`.

However, in the examples above that index the 2-dimensional `mtcars` object, the square brackets are used to index **row** and **column** positions. The syntax for this is `[rows, columns]`. We will be using such indexing throughout this module.

You can ask R to return you specific rows and columns by different ways:

```
mtcars[c(2,9), 3:7]
```

| | disp | hp | drat | wt | qsec |
|---------------|-------|-----|------|-------|-------|
| Mazda RX4 Wag | 160.0 | 110 | 3.90 | 2.875 | 17.02 |
| Merc 230 | 140.8 | 95 | 3.92 | 3.150 | 22.90 |

```
mtcars[3:6, c("disp","hp","qsec")]
```

| | disp | hp | qsec |
|-------------------|------|-----|-------|
| Datsun 710 | 108 | 93 | 18.61 |
| Hornet 4 Drive | 258 | 110 | 19.44 |
| Hornet Sportabout | 360 | 175 | 17.02 |
| Valiant | 225 | 105 | 20.22 |

```
mtcars[, c("wt","gear","cyl")]
```

| | wt | gear | cyl |
|---------------------|-------|------|-----|
| Mazda RX4 | 2.620 | 4 | 6 |
| Mazda RX4 Wag | 2.875 | 4 | 6 |
| Datsun 710 | 2.320 | 4 | 4 |
| Hornet 4 Drive | 3.215 | 3 | 6 |
| Hornet Sportabout | 3.440 | 3 | 8 |
| Valiant | 3.460 | 3 | 6 |
| Duster 360 | 3.570 | 3 | 8 |
| Merc 240D | 3.190 | 4 | 4 |
| Merc 230 | 3.150 | 4 | 4 |
| Merc 280 | 3.440 | 4 | 6 |
| Merc 280C | 3.440 | 4 | 6 |
| Merc 450SE | 4.070 | 3 | 8 |
| Merc 450SL | 3.730 | 3 | 8 |
| Merc 450SLC | 3.780 | 3 | 8 |
| Cadillac Fleetwood | 5.250 | 3 | 8 |
| Lincoln Continental | 5.424 | 3 | 8 |
| Chrysler Imperial | 5.345 | 3 | 8 |
| Fiat 128 | 2.200 | 4 | 4 |
| Honda Civic | 1.615 | 4 | 4 |
| Toyota Corolla | 1.835 | 4 | 4 |
| Toyota Corona | 2.465 | 3 | 4 |
| Dodge Challenger | 3.520 | 3 | 8 |
| AMC Javelin | 3.435 | 3 | 8 |
| Camaro Z28 | 3.840 | 3 | 8 |
| Pontiac Firebird | 3.845 | 3 | 8 |
| Fiat X1-9 | 1.935 | 4 | 4 |

| | | | |
|----------------|-------|---|---|
| Porsche 914-2 | 2.140 | 5 | 4 |
| Lotus Europa | 1.513 | 5 | 4 |
| Ford Pantera L | 3.170 | 5 | 8 |
| Ferrari Dino | 2.770 | 5 | 6 |
| Maserati Bora | 3.570 | 5 | 8 |
| Volvo 142E | 2.780 | 4 | 4 |

1.2.5.7 Packages

The **base** installation of R includes many functions and commands. However, more often we are interested in using some particular functionality, encoded into **packages** contributed by the R developer community. Installing packages for the first time can be done at the command line in the R console using the `install.packages` command as in the example below to install the **tmmap** library or via the RStudio menu via **Tools > Install Packages**.

When you install these packages it is strongly suggested you also install the *dependencies*. These are other packages that are required by the package that is being installed. This can be done by selecting check the box in the menu or including `dep=TRUE` in the command line as below (don't run this yet!):

```
# don't run this!
install.packages("tidyverse", dep = TRUE)
```

You may have to set a **mirror** site from which the packages will be downloaded to your computer. Generally you should pick one that is nearby to you.

Further descriptions of packages, their installation and their data structures will be given as needed in the practicals. There are literally 1000s of packages that have been contributed to the R project by various researchers and organisations. These can be located by name at http://cran.r-project.org/web/packages/available_packages_by_name.html if you know the package you wish to use. It is also possible to search the CRAN website to find packages to perform particular tasks at <http://www.r-project.org/search.html>. Additionally many packages include user guides and vignettes as well as a PDF document describing the package and listed at the top of the index page of the help files for the package.

As well as **tidyverse** you should install the **sf** package and dependencies. So we have 2 packages to install:

- **sf** for spatial data and spatial objects
- **tidyverse** for lots of lovely data science things - see <https://www.tidyverse.org>

You could do this in one go and this will take a bit of time:


```
install.packages(c("sf", "tidyverse"), dep = TRUE)
```

Remember: you will only have to install a package once!! So when the above code has run in your script you should comment it out. For example you might want to include something like the below in your R script.

```
# packages only need to be loaded once  
# install.packages(c("sf", "tidyverse"), dep = TRUE)
```

Once the package has been installed on your computer then the package can be called using the `library()` function into each of your R sessions as below.

```
library(tidyverse)  
library(sf)
```

1.3 Knowing Merseyside

1.3.1 Merseyside districts

Now we use these basic R command and newly installed packages to start our initial exploration by using some existing secondary dataset from the Census 2021.

In R we normally read in tabular dataset from .csv format. In your [ENVS162 Canvas page](#) find Week 1 -> Practical 1 Dataset, download the four datasets to your current working folder on your M drive (ENVS162 - Week 1). You may first identify one .csv dataset: **merseyside.csv**. You can open them in excel to have a look, but here we are using R instead of Excel to load and examine them.

1.3.1.1 Loading tabular data

The survey data can be loaded into RStudio using the `read.csv` function.

However, you will need to tell R where to get the data from. The easiest way to do this is to use the menu if the R script file is open. Go to **Session > Set Working Directory > To Source File Location** to set the working directory to the location where your `week1.R` script is saved. When you do this you will see line of code print out in the Console (bottom left pane) similar to `setwd("SomeFilePath")`. You can copy this line of code to your script and paste into the line above the line calling the `read.csv` function.

```
# use read.csv to load a CSV file
# this is assignment to an object called `df`
df = read.csv(file = "merseyside.csv", stringsAsFactors = TRUE)
```

The `stringsAsFactors = TRUE` parameter tells R to read any character or text variables as classes or categories and not as just text.

You could inspect the help for the `read.csv` function to see the different parameters and their default values:

```
help(read.csv)
# or
?read.csv
```

Functions always return something and in this case `read.csv()` function has returned a tabular R object with 5 records and 12 fields. This has been *assigned to df*.

Finally in this section, let's have a look at the data. This can be done in a number of ways.

- you could look at the `df` object by entering `df` in the Console. However this is not particularly helpful as it simply prints out everything that is in `df` to the Console.
- you could click on the `df` object in the Environment pane and this shows the structure of the attributes in different fields.
- you could click on the little grid-like icon next `df` in the Environment pane to get a **View** of the data and remember to close the tab that opens!.
- or you could use some code as in the examples below.

First, let's have a look at the internal structure of the data using the `str` function:

```
str(df)
```

```
'data.frame':  5 obs. of  12 variables:
 $ LAD21CD      : Factor w/ 5 levels "E08000011","E08000012",...: 1 2 3 4 5
 $ District     : Factor w/ 5 levels "Knowsley","Liverpool",...: 1 2 4 3 5
 $ Population   : int  154519 486089 183248 279234 320196
 $ Households   : int  66073 207491 81011 123075 143253
 $ Working_population: int  69495 205749 82622 124596 139500
 $ Students     : int  7050 59628 7582 12636 14642
 $ Unemployed   : int  3852 13894 4076 6143 6542
 $ Age_over_65  : int  26242 74322 37642 64763 70391
 $ Disability   : int  34990 105962 40829 61134 73088
```

```
$ No_central_heating: int 1020 4822 1003 1965 2125
$ Overcrowding      : int 1892 7352 1888 2700 2355
$ Working_from_home : int 14880 53721 18973 34750 37299
```

There is other ways to get info about the number of rows and columns:

```
nrow(df)
```

```
[1] 5
```

```
ncol(df)
```

```
[1] 12
```

```
#or both row and col
dim(df)
```

```
[1] 5 12
```

The `head` function does this by printing out the first six records of the data table and you may need to scroll up and down in the Console pane to see all of what is returned.

```
head(df)
```

| | LAD21CD | District | Population | Households | Working_population | Students |
|---|-------------------|-------------|------------|--------------------|--------------------|----------|
| 1 | E08000011 | Knowsley | 154519 | 66073 | 69495 | 7050 |
| 2 | E08000012 | Liverpool | 486089 | 207491 | 205749 | 59628 |
| 3 | E08000013 | St. Helens | 183248 | 81011 | 82622 | 7582 |
| 4 | E08000014 | Sefton | 279234 | 123075 | 124596 | 12636 |
| 5 | E08000015 | Wirral | 320196 | 143253 | 139500 | 14642 |
| | Unemployed | Age_over_65 | Disability | No_central_heating | Overcrowding | |
| 1 | 3852 | 26242 | 34990 | 1020 | 1892 | |
| 2 | 13894 | 74322 | 105962 | 4822 | 7352 | |
| 3 | 4076 | 37642 | 40829 | 1003 | 1888 | |
| 4 | 6143 | 64763 | 61134 | 1965 | 2700 | |
| 5 | 6542 | 70391 | 73088 | 2125 | 2355 | |
| | Working_from_home | | | | | |
| 1 | 14880 | | | | | |
| 2 | 53721 | | | | | |
| 3 | 18973 | | | | | |
| 4 | 34750 | | | | | |
| 5 | 37299 | | | | | |

Another way to explore the data is through the `summary` function:

```
summary(df)
```

| LAD21CD | District | Population | Households |
|--------------------|--------------------|----------------|-------------------|
| E08000011:1 | Knowsley :1 | Min. :154519 | Min. : 66073 |
| E08000012:1 | Liverpool :1 | 1st Qu.:183248 | 1st Qu.: 81011 |
| E08000013:1 | Sefton :1 | Median :279234 | Median :123075 |
| E08000014:1 | St. Helens:1 | Mean :284657 | Mean :124181 |
| E08000015:1 | Wirral :1 | 3rd Qu.:320196 | 3rd Qu.:143253 |
| | | Max. :486089 | Max. :207491 |
| Working_population | Students | Unemployed | Age_over_65 |
| Min. : 69495 | Min. : 7050 | Min. : 3852 | Min. :26242 |
| 1st Qu.: 82622 | 1st Qu.: 7582 | 1st Qu.: 4076 | 1st Qu.:37642 |
| Median :124596 | Median :12636 | Median : 6143 | Median :64763 |
| Mean :124392 | Mean :20308 | Mean : 6901 | Mean :54672 |
| 3rd Qu.:139500 | 3rd Qu.:14642 | 3rd Qu.: 6542 | 3rd Qu.:70391 |
| Max. :205749 | Max. :59628 | Max. :13894 | Max. :74322 |
| Disability | No_central_heating | Overcrowding | Working_from_home |
| Min. : 34990 | Min. :1003 | Min. :1888 | Min. :14880 |
| 1st Qu.: 40829 | 1st Qu.:1020 | 1st Qu.:1892 | 1st Qu.:18973 |
| Median : 61134 | Median :1965 | Median :2355 | Median :34750 |
| Mean : 63201 | Mean :2187 | Mean :3237 | Mean :31925 |
| 3rd Qu.: 73088 | 3rd Qu.:2125 | 3rd Qu.:2700 | 3rd Qu.:37299 |
| Max. :105962 | Max. :4822 | Max. :7352 | Max. :53721 |

Finally in this section, we come back to the dollar sign (\$). This is used to refer to or *extract* an individual named field or variable in an R object, like `df`.

The code below prints out the Population attribute and generates a summary of its values:

```
# extract an individual variable
df$Population
```

```
[1] 154519 486089 183248 279234 320196
```

```
# generate a summary of an individual variable
summary(df$Population)
```

| | | | | | |
|--------|---------|--------|--------|---------|--------|
| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
| 154519 | 183248 | 279234 | 284657 | 320196 | 486089 |

And of course we can use such operations to *assign* the result to new R objects. The code below extracts three variables from `df`, assigns them to `x`, `y` and `z`, and then uses the `data.frame` function to convert these into a new `data.frame` object called `my_df`

```
# extract three variables, assigning them to temporary R objects
x = df$District
y = df$Working_population
z = df$Students
# create a data.frame from these, naming the new variables
my_df = data.frame(district = x, worker = y, student = z)
```

You should have a look at what you have created:

```
head(my_df)
```

| | district | worker | student |
|---|------------|--------|---------|
| 1 | Knowsley | 69495 | 7050 |
| 2 | Liverpool | 205749 | 59628 |
| 3 | St. Helens | 82622 | 7582 |
| 4 | Sefton | 124596 | 12636 |
| 5 | Wirral | 139500 | 14642 |

```
summary(my_df)
```

| | district | worker | student |
|------------|----------|----------------|---------------|
| Knowsley | :1 | Min. : 69495 | Min. : 7050 |
| Liverpool | :1 | 1st Qu.: 82622 | 1st Qu.: 7582 |
| Sefton | :1 | Median :124596 | Median :12636 |
| St. Helens | :1 | Mean :124392 | Mean :20308 |
| Wirral | :1 | 3rd Qu.:139500 | 3rd Qu.:14642 |
| | | Max. :205749 | Max. :59628 |

The temporary R objects can be removed from the Environment using the `rm` function and a *combine* vector function, `c()` that you encountered in Week 19, that takes a vector of object names (hence they are in quotes) as its arguments.

```
rm(list = c("x", "y", "z"))
```

1.3.1.2 Basic data manipulation

Now we can do some basic data manipulation to know Merseyside more from the data perspective.

What is the total population in Merseyside?

```
sum(df$Population)
```

```
[1] 1423286
```

What is the total number of full-time students in Merseyside?

```
sum(df$Students)
```

```
[1] 101538
```

Then, we can calculate the total number of workers that working from home:

```
sum(df$Working_from_home)
```

```
[1] 159623
```

What is the proportion of working population actually work from home in Merseyside? Yes, we need to use a division calculation of the total number of working from home vs. all the working population. R can do it by:

```
sum(df$Working_from_home) / sum(df$Working_population)
```

```
[1] 0.2566443
```

So the answer is 25.7% for the whole Merseyside - but which district has the highest proportion and which as the lowest? You may have your own guessing. But let R do the calculation:

```
df$Prop.WFH = df$Working_from_home / df$Working_population * 100 #add a new column called Pro  
df #print out the df
```

| | LAD21CD | District | Population | Households | Working_population | Students |
|---|-----------|------------|------------|------------|--------------------|----------|
| 1 | E08000011 | Knowsley | 154519 | 66073 | 69495 | 7050 |
| 2 | E08000012 | Liverpool | 486089 | 207491 | 205749 | 59628 |
| 3 | E08000013 | St. Helens | 183248 | 81011 | 82622 | 7582 |
| 4 | E08000014 | Sefton | 279234 | 123075 | 124596 | 12636 |
| 5 | E08000015 | Wirral | 320196 | 143253 | 139500 | 14642 |

| | Unemployed | Age_over_65 | Disability | No_central_heating | Overcrowding |
|---|------------|-------------|------------|--------------------|--------------|
| 1 | 3852 | 26242 | 34990 | 1020 | 1892 |
| 2 | 13894 | 74322 | 105962 | 4822 | 7352 |
| 3 | 4076 | 37642 | 40829 | 1003 | 1888 |
| 4 | 6143 | 64763 | 61134 | 1965 | 2700 |
| 5 | 6542 | 70391 | 73088 | 2125 | 2355 |

| | Working_from_home | Prop.WFH |
|---|-------------------|----------|
| 1 | 14880 | 21.41161 |
| 2 | 53721 | 26.10997 |
| 3 | 18973 | 22.96362 |
| 4 | 34750 | 27.89014 |
| 5 | 37299 | 26.73763 |

Here we ask R to add a new column named `Prop.WFH` which is the working from home proportion that calculated by the number of working from home people in each district divided by the total working population in that district. The `* 100` convert the rate in the percentage number. R will automatically do it row-by-row. We then print out the `df`, you may find at the very right end of the tabular, there is a new column called `Prop.WFH`.

For a very small dataframe like this, we can also using `View()` to open a new tab to review the data, where each column can be sorted from largest to smallest or vice versa. Try viewing it and find the newly created column `Prop.WFH`. Click on the column name, you should see it is sorted from highest to lowest, and click again, the ranking is reversed.

```
View(df)
```

1.3.1.3 Your first map for Merseyside

Now let's try to do our first map in R and allow yourself know more about Merseyside.

We will use the library `sf` and `tmap` to help us at here. Run the install codes if you haven't install them. Remember: you will only have to install a package once!!

```
install.packages("tmap", dep = TRUE)
```

Check the package version of `tmap`, as here we need to use `tmap` over 4.0 version.

```
packageVersion("tmap") # the version should over 4.0
```

```
[1] '4.1'
```

When they have been installed, we can start using them

```
library(tidyverse)
library(sf)
library(tmap)
```

You may find in Week 1 data, we have another file named *merseyside_districts.gpkg*. A GeoPackage (GPKG) is a file-based format designed for storing geographic data. It supports the efficient storage and exchange of spatial datasets and can be readily used across GIS software such as QGIS and ArcGIS, as well as in programming environments including R and Python.

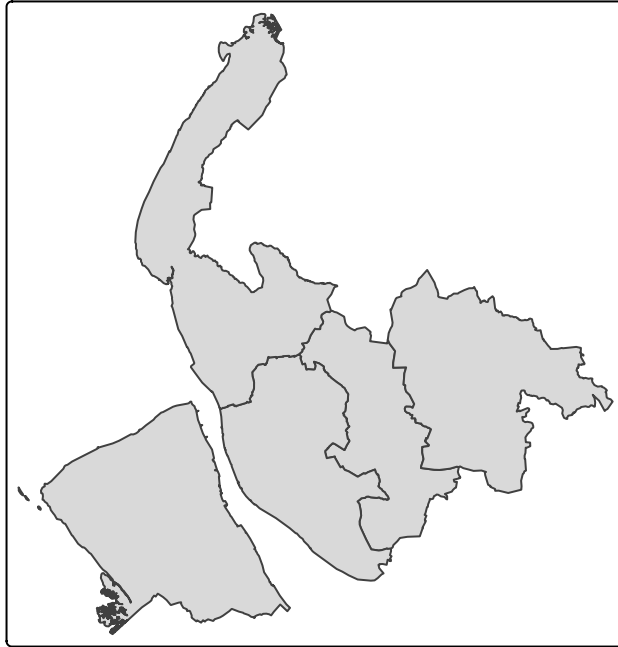
We first read it in by using the `st_read()` command in library `sf`.

```
sf <- st_read("merseyside_districts.gpkg")
```

```
Reading layer `lad_may_2025_uk_bgc_v2_4306843991635065087__lad_may_2025_uk_bgc_v2' from data
  using driver `GPKG'
Simple feature collection with 5 features and 8 fields
Geometry type: MULTIPOLYGON
Dimension:      XY
Bounding box:   xmin: 318351.7 ymin: 377515.4 xmax: 361796.3 ymax: 422866.5
Projected CRS:  OSGB36 / British National Grid
```

The fastest way to map it is the `qtm()` function.

```
qtm(sf)
```

You can also add the district names on the map - which column in the sf contains district name? Use `names(sf)` to check for it.

Yes, the column should be LAD25NM. Now let's ask `qtm()` to also show the district names.

```
qtm(sf, text="LAD25NM")
```



But what if we want to make some meaningful maps, rather than just the boundaries of these five districts of Merseyside?

1.3.1.4 Link tabular data to geographical boundaries

Recall that in our `df`, we have 14 columns, containing different information about the districts. We can get all their names by using `names()`.

```
names(df)
```

```
[1] "LAD21CD"          "District"          "Population"
[4] "Households"       "Working_population" "Students"
[7] "Unemployed"       "Age_over_65"       "Disability"
[10] "No_central_heating" "Overcrowding"      "Working_from_home"
[13] "Prop.WFH"
```

We can do the same thing for our geographical dataset `sf` to see what it includes:

```
names(sf)
```

```
[1] "LAD25CD" "LAD25NM" "LAD25NMW" "BNG_E" "BNG_N" "LONG" "LAT"
[8] "GlobalID" "geom"
```

We can also show the whole `sf` as

```
sf
```

Simple feature collection with 5 features and 8 fields

Geometry type: MULTIPOLYGON

Dimension: XY

Bounding box: xmin: 318351.7 ymin: 377515.4 xmax: 361796.3 ymax: 422866.5

Projected CRS: OSGB36 / British National Grid

| | LAD25CD | LAD25NM | LAD25NMW | BNG_E | BNG_N | LONG | LAT |
|---|-----------|------------|----------|--------|--------|-----------|----------|
| 1 | E08000011 | Knowsley | | 344762 | 393778 | -2.832979 | 53.43789 |
| 2 | E08000012 | Liverpool | | 339359 | 390556 | -2.913680 | 53.40833 |
| 3 | E08000013 | St. Helens | | 353413 | 395992 | -2.703093 | 53.45862 |
| 4 | E08000014 | Sefton | | 334282 | 398835 | -2.991771 | 53.48213 |
| 5 | E08000015 | Wirral | | 329109 | 386965 | -3.067034 | 53.37478 |

| | GlobalID | geom |
|---|--|--------------------------------|
| 1 | {B4196BFE-EE90-4C31-ABD5-C7E743AE2F9B} | MULTIPOLYGON (((341447.1 40... |
| 2 | {4FB47E7A-EF4E-4B9E-BF75-D4FC059CDE61} | MULTIPOLYGON (((338860.9 39... |
| 3 | {943F0C6B-EB30-4C00-A42B-F6B3AEC3EFEE} | MULTIPOLYGON (((349111.4 40... |
| 4 | {C6FD073B-CBEB-4E78-934A-A8FD11A20FOA} | MULTIPOLYGON (((336374.5 42... |
| 5 | {88E9328B-371C-469C-91F1-3479C77D6950} | MULTIPOLYGON (((331364.9 39... |

Now we see that `sf` includes also the five districts, but also other geographical information. You may notice that although different column names, the first two columns of both `df` and `sf` are the district code and district name. This means what potentially we can link this two dataset together - appendix the `df` to `sf` to enrich the attributes of our geographical dataset.

```
merseyside <- left_join(sf, df, by=c("LAD25NM"="District"))
```

let's check out the new `sf2` by `View()` it:

```
View(merseyside)
```

In the open tab, we see all the `df` columns are now also attached to the `sf`, linking by the district names.

1.3.1.5 Choropleth map of Merseyside districts

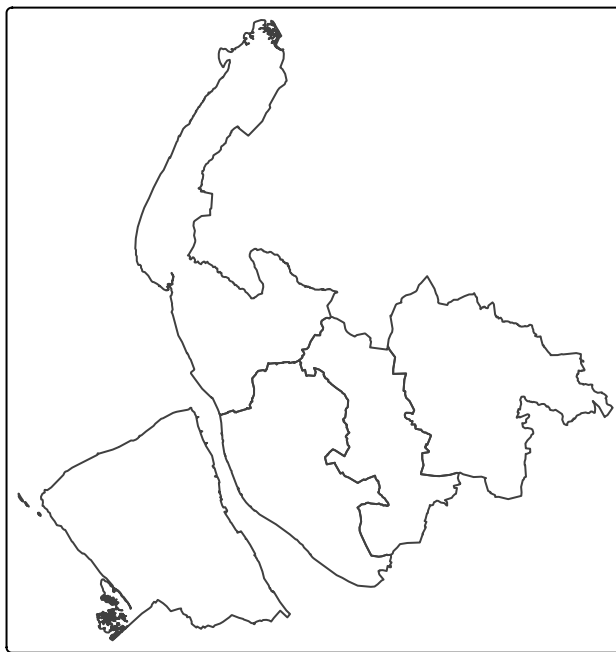
Now, we can use those new columns we attached from `df` to `sf2` to make some meaningful choropleth maps! Here we make use of the mapping functions in `tmap` (Remember to run `library(tmap)` if you haven't) to do the work for us.

`tmap` has a basic syntax (again, do not run this code - its is simply showing the syntax of `tmap`):

```
# don't run this or write this into your script!  
tm_shape(data = <data>)+  
  tm_<function>(<variable to be mapped>)
```

For example, to map the boundaries of `merseyside`:

```
tm_shape(merseyside) +  
  tm_borders()
```



To add label of district:

```
tm_shape(merseyside) +  
  tm_borders() +  
  tm_text("LAD25NM")
```

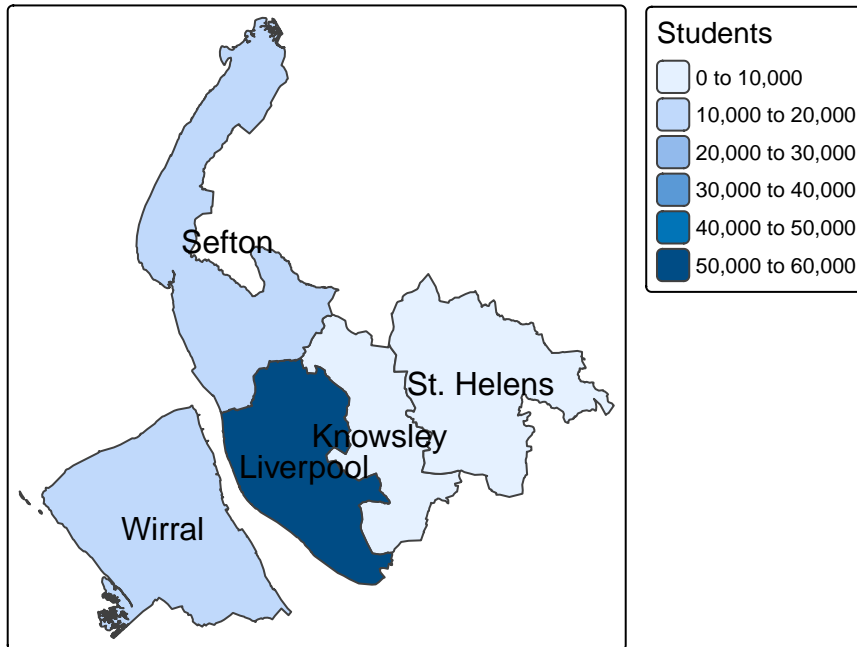


You might assume the quick mapping function `qtm()` can achieve the same result, but `tmap` provides far more flexibility when it comes to aesthetic customization. The easiest way to illustrate `tmap` is through some examples.

Let's start with a simple choropleth map, by using `tmap` to show the distribution of a continuous variable in different elements of the spatial data (here are the data Merseyside districts are polygons).

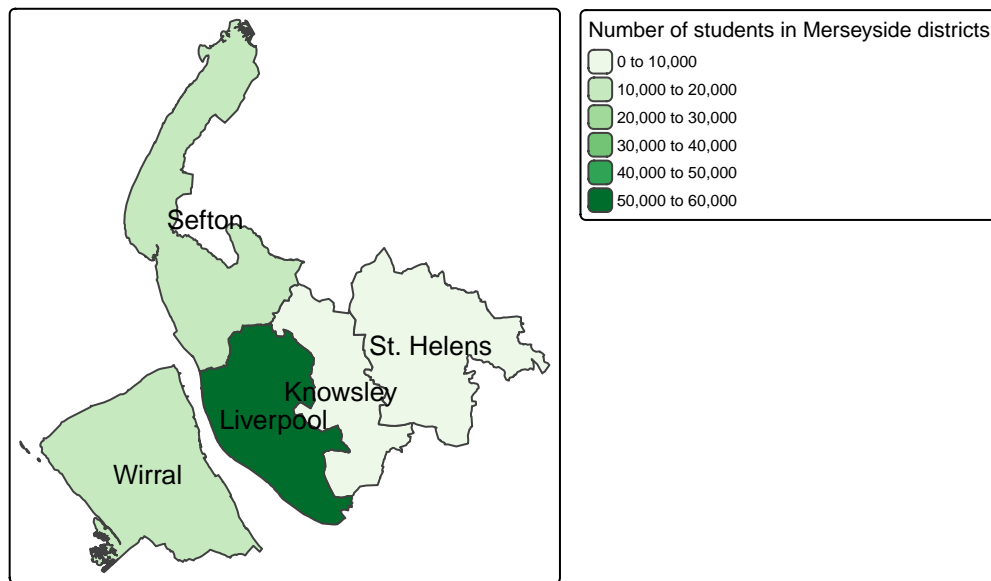
The code below maps 'Students' as in the Merseyside districts, and shows the district names of each polygon from 'LAD25NM' columns. The map below indicates that Liverpool has the highest number of full-time students while Knowsley and St.Helens have the least.

```
tm_shape(merseyside) +  
  tm_polygons(fill = "Students") + # Variable to map  
  tm_text("LAD25NM")              # Variable to label
```



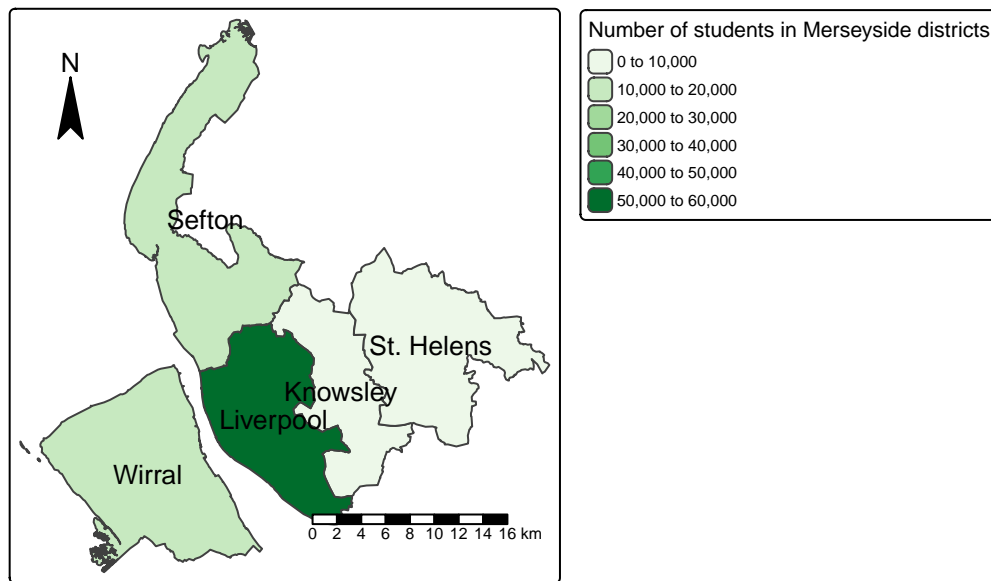
By default tmap picks a shading scheme, the class breaks and places a legend somewhere. All of these can be changed. The code below allocates the tmap plot to **map1** (Map 1), change the legend title as “Number of students in Merseyside districts”, and then prints it:

```
map1 = tm_shape(merseyside) +
  tm_polygons(fill="Students",
              fill.scale = tm_scale(values = "Greens"), # change palette to greens
              fill.legend = tm_legend(title = "Number of students in Merseyside districts")
              ) + # Legend title
  tm_text("LAD25NM",size=0.8) #size down the label slightly
map1
```



And of course many other elements included either by running the code snippet defining `map1` above with additional lines or by simply adding them as in the code below:

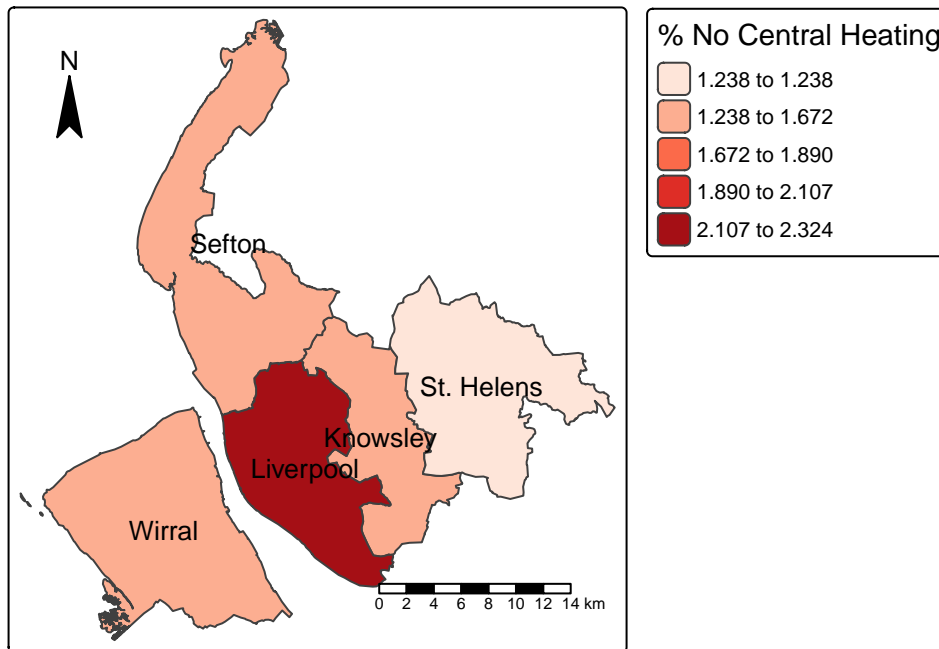
```
map1 +  
  tm_scalebar(position = c("right", "bottom")) +  
  tm_compass(position = c("left", "top")) # Use "top", "center", or "bottom"
```



We can also create new variable to the dataset and then map it. The below code chunk first creates a new column, “NoCentralHeating_rate”, by dividing the number of households without access to central heating by the total number of households in each district; it then uses `tmap` to make a map of the proportion of households without central heating across districts in Merseyside:

```
merseyside$NoCentralHeating_rate = merseyside$No_central_heating / merseyside$Households * 100

map2 = tm_shape(merseyside) +
  tm_polygons(fill="NoCentralHeating_rate",
              fill.scale = tm_scale(values = "Reds", style = "jenks"), #use jenks classification
              fill.legend = tm_legend(title = "% No Central Heating")) +
  tm_text("LAD25NM",size=0.8) +
  tm_scalebar(position = c("right", "bottom")) + # Add a scale bar at the top-right corner
  tm_compass(position = c("left", "top")) # Add a compass rose at the top-right corner
map2
```

1.3.2 Merseyside neighbourhoods

Now let's read in the neighbourhood-level datasets, which include a .csv file of local statistics and the corresponding geographical boundaries.

```
lsoa_df <- read.csv("merseyside_lsoa.csv")
lsoa_sf <- st_read("LSOA_boundaries.gpkg")
```

```
Reading layer `merseyside_LSOA' from data source
  `C:\Users\ziye\Documents\GitHub\quant\labs\LSOA_boundaries.gpkg'
  using driver `GPKG'
Simple feature collection with 923 features and 4 fields
Geometry type: MULTIPOLYGON
Dimension:      XY
Bounding box:   xmin: -3.200368 ymin: 53.2963 xmax: -2.576743 ymax: 53.6982
Geodetic CRS:   WGS 84
```

First, we take a look at the .csv dataset, which as been read into R as lsoa_df:

```
View(lsoa_df)
```

or check the structure of the dataset:

```
str(lsoa_df)
```

```
'data.frame':  923 obs. of  11 variables:
 $ LSOA21CD      : chr  "E01006416" "E01006418" "E01006434" "E01006435" ...
 $ Population    : int  1520 1315 1519 1524 1150 1654 1450 1581 1421 1373 ...
 $ Households    : int  678 567 652 663 490 695 592 622 809 618 ...
 $ Working_population: int  588 547 660 581 546 766 558 570 524 481 ...
 $ Students      : int  59 64 69 82 51 66 65 89 52 72 ...
 $ Unemployed    : num  3.57 2.74 5.11 3.43 1.62 ...
 $ Age_over_65   : num  14.8 17.5 11.7 19.9 26.3 ...
 $ Disability     : num  27.1 30 23.4 29 20.5 ...
 $ No_central_heating: int  16 14 16 7 3 8 9 9 12 9 ...
 $ Overcrowding  : int  24 24 35 28 13 21 29 23 31 22 ...
 $ Working_from_home : int  91 84 102 96 165 171 101 64 66 48 ...
```

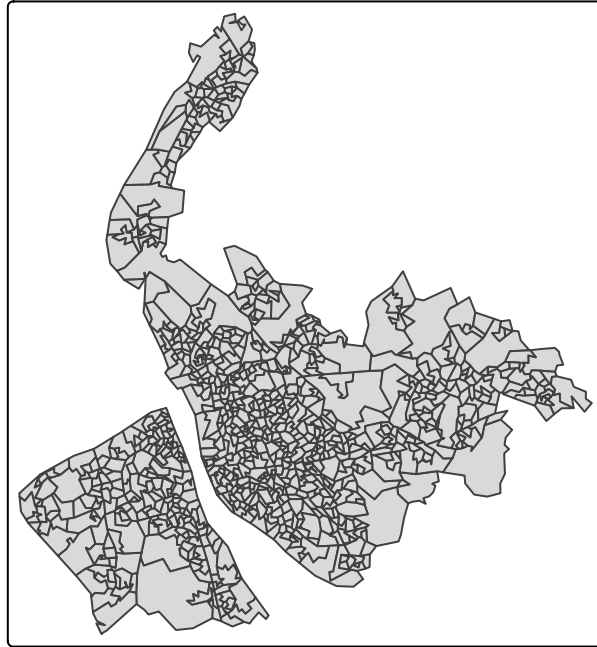
So now, you know how many LSOAs in Merseyside? Yes, there are 923 LSOAs. As we introduced in the Week 1 lecture, LSOA means Super Output Area Lower Area and is commonly used in the Census statistics. Each LSOA represents 1,000 to 3,000 people or 400 to 1,200 households in England and Wales.

```
dim(lsoa_df)
```

```
[1] 923  11
```

Use the quick mapping function `qtm()` to quickly inspect the geographical boundary dataset `lsoa_sf`.

```
qtm(lsoa_sf)
```



check how many LSOAs in the boundary dataset - there should also be 923.

```
nrow(lsoa_sf)
```

```
[1] 923
```

To familiarise yourself with the structures of both datasets, we can use the `names()` command

```
names(lsoa_df)
```

```
[1] "LSOA21CD"          "Population"        "Households"
[4] "Working_population" "Students"          "Unemployed"
[7] "Age_over_65"       "Disability"        "No_central_heating"
[10] "Overcrowding"      "Working_from_home"
```

```
names(lsoa_sf)
```

```
[1] "LSOA21CD" "LSOA21NM" "LAD23CD"  "LAD23NM"  "geom"
```

You may find that both dataset are recorded at the LSOA level, with LSOA21CD as the key column. As we did with the district-level dataset, we can use `left_join()` to join these two dataset by their sharing field - LSOA21CD:

```
lsoa <- left_join(lsoa_sf, lsoa_df, by="LSOA21CD")
```

Now let's check the columns of new dataframe `lsoa`:

```
names(lsoa)
```

```
[1] "LSOA21CD"          "LSOA21NM"          "LAD23CD"
[4] "LAD23NM"           "Population"         "Households"
[7] "Working_population" "Students"           "Unemployed"
[10] "Age_over_65"       "Disability"         "No_central_heating"
[13] "Overcrowding"      "Working_from_home"  "geom"
```

Or open a new tab to view the newly created dataset `lsoa` by

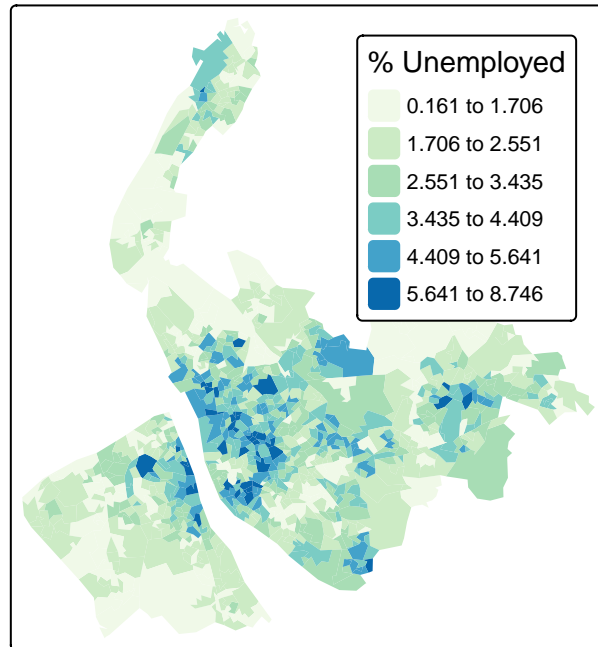
```
View(lsoa)
```

We can see that some columns contain counts, such as the number of residential population, number of households, number of working population, and number of students. Other columns are expressed as percentages, including unemployment, population aged 65 and over, disability, households without central heating, overcrowded households, and people working from home.

1.3.2.1 Making maps across LSOAs in Merseyside

Using the `Unemployed` column, we can create a map of the unemployment rate across neighbourhoods in Merseyside. Instead of using the default equal-interval breaks, this time we will use a jenks classification with six categories.

```
map3 = tm_shape(lsoa) +
  tm_fill(
    fill = "Unemployed",
    fill.scale = tm_scale(values = "GnBu",
                          style = "jenks",
                          n = 6), #use jenks classification of 6 categories
    fill.legend = tm_legend(title = "% Unemployed")
  ) +
  tm_layout(legend.position = c("right", "top"))
map3
```



The above code uses `tm_layout(legend.position = c("right", "top"))` to move the legend inside the map frame, positioning it at the right-top corner.

Replace `tm_fill()` to `tm_polygons()` to see how the map changes?

`tm_polygons()` is a condensed version of `tm_fill()` + `tm_border()`. Here if you want to show all the LSOA borders, use `tm_polygons()` instead of `tm_fill()`.

1.3.2.2 Overlapping tmap objects

`tmap` also supports adding or overlaying other data, such as boundaries. Because these are additional spatial data layers, they need to be added with `tm_shape()` followed by the usual function.

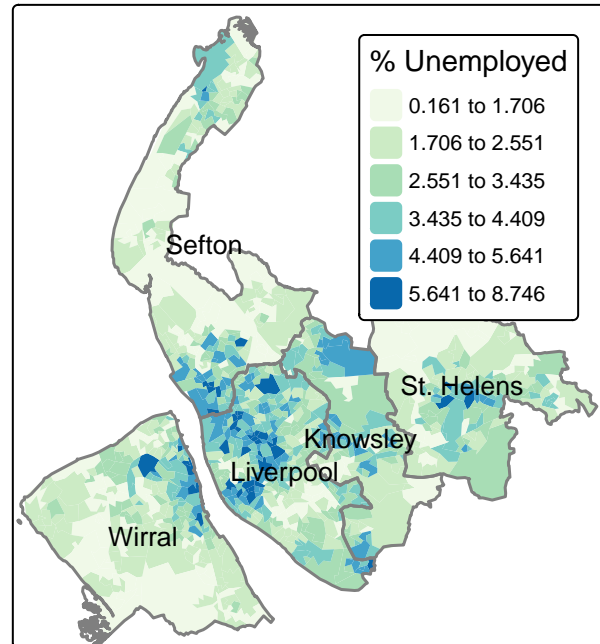
Remember we use the code chunk to make the district boundaries of Merseyside? This time let's change the aesthetic by using grey color as the border color and increase the line width, then we save it also as a `tmap` object called `map_district`:

```
map_district = tm_shape(merseyside) +
  tm_borders(col = "grey50", lwd=1.5) + #border color as grey, line width as 1.5
  tm_text("LAD25NM", size = 0.8)
map_district
```



To display both `tmap` layers together, we can proceed as follows:

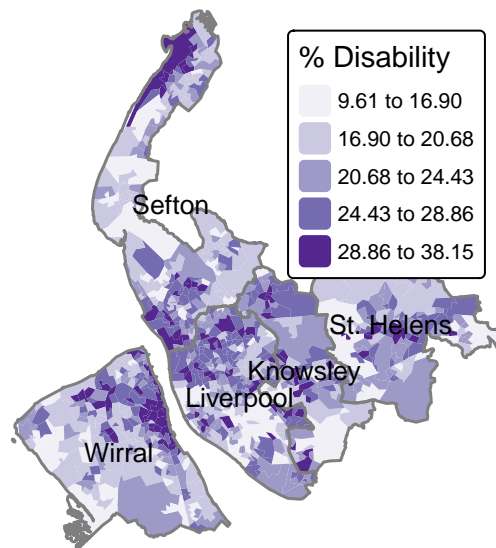
```
map3 + map_district
```



Now, let's move to making some more maps - this time showing the proportion of disability across neighbourhoods in Merseyside. Referring back to the columns in `lsoa`, this time we use the `Disability` variable. The code below applies a jenks classification and use a different color palette `Purples`. Also the map frame is removed by `frame = FALSE`.

```
tm_shape(lsoa) +  
  tm_fill(fill = "Disability",  
          fill.scale = tm_scale(values="Purples",  
                                style = "jenks",  
                                n=5),  
          fill.legend = tm_legend(title = "% Disability")  
        ) +  
  tm_layout(main.title = "Merseyside",#add a main title  
            legend.position = c("right", "top"),  
            frame = FALSE)+  
  map_district
```

Merseyside



1.3.2.3 Create new variables to make maps

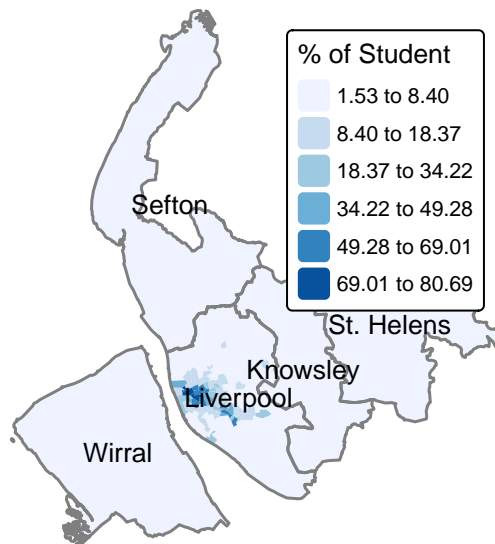
Similarly, we can make maps from new columns we made ourselves. For example, we can calculate the percentage of student by adding a new column to the dataframe:

```
lsoa$student.perc = lsoa$Students / lsoa$Population * 100
```

To make a map to visualisation the spatial distribution of student percentage. The code below uses `n=6` to increase the classification categories to 6 rather than default 5.

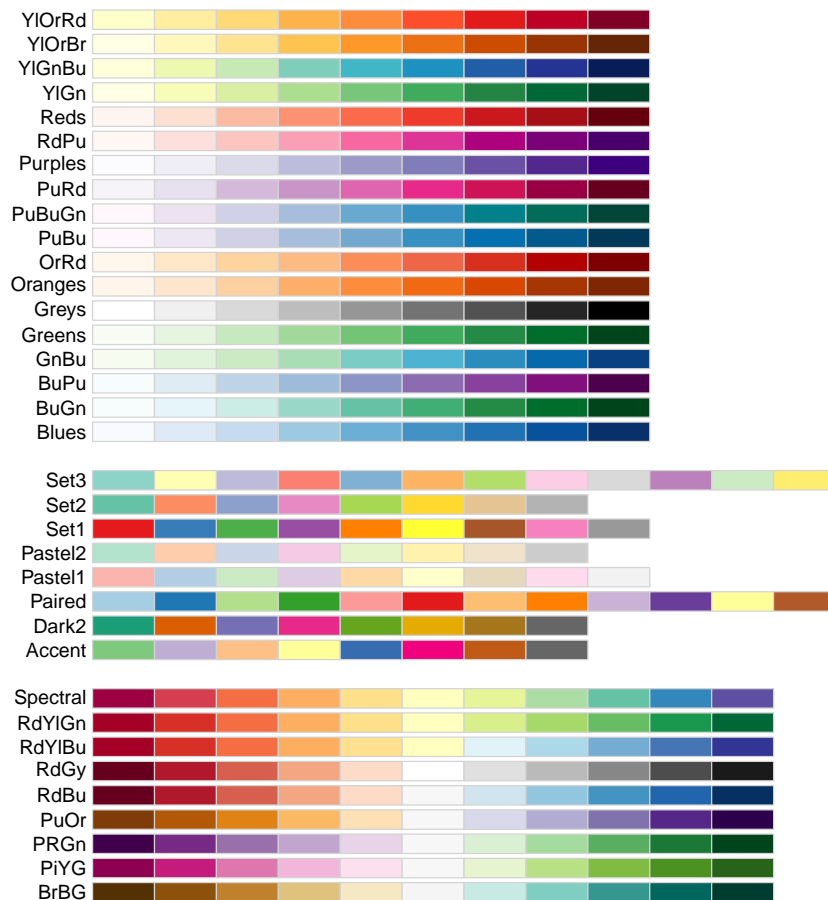
```
tm_shape(lsoa) +  
  tm_fill(fill = "student.perc",  
          fill.scale = tm_scale(values="Blues",  
                                style = "jenks",  
                                n=6),  
          fill.legend = tm_legend(title = "% of Student",  
                                  title.size = 0.8) #legend title change smaller font  
  ) +  
  tm_layout(main.title = "Merseyside",  
            frame = FALSE,  
            legend.position = c("right", "top"))+  
  map_district
```

Merseyside



To change the palette, RColorBrewer provides different palette choices:

```
RColorBrewer::display.brewer.all()
```

1.3.2.4 In a nutshell

If we want to make a map to show the rate of no central heating households in all the neighbourhoods in Merseyside, we need to create a new variable `no.central.heating.perc` as the result of dividing households without central heating by total households in each LSOA. The code below combines all the cartographic elements together:

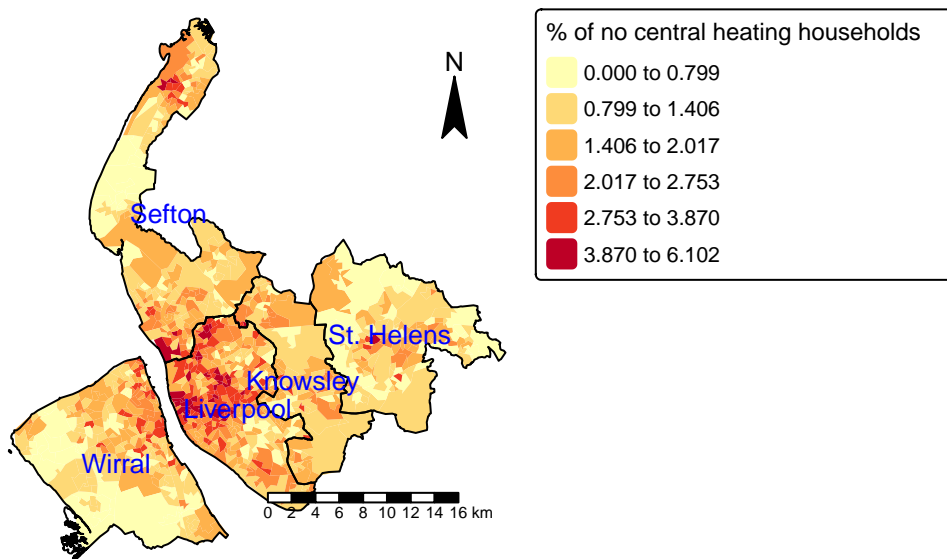
```

lsoa$no.central.heating.perc = lsoa$No_central_heating / lsoa$Households * 100

tm_shape(lsoa) +
  tm_fill(fill = "no.central.heating.perc",
          fill.scale = tm_scale(values="YlOrRd",
                                style = "jenks",n=6),
          fill.legend = tm_legend(title = "% of no central heating households",
                                  title.size = 0.8)
        ) +
  tm_layout(main.title = "Merseyside",
            main.title.size=1.2,
            frame = FALSE) +
  tm_compass(position = c("right", "top")) +
  tm_scalebar(position = c("right", "bottom")) +
  tm_shape(merseyside) + # Add another spatial layer (Merseyside boundary)
  tm_borders(col = "black", lwd = 1) + # Draw the boundaries with black lines of width 1
  tm_text("LAD25NM",col = "blue",size = 0.8)

```

Merseyside



1.4 Summary

The aim of this session has been to familiarise you with the R environment if you have not used R before. If you have but not for a while, then hopefully this has acted as a refresher. Some key things to take away are:

- R is a learning curve, and like driving the more your practice the better you become.
- Your job is to try to **understand** what the code is doing and **not** to remember the code.
- To help with this, you should add your own comments to the script to help you understand what is going on when you return to them. Comments are prefaced by a hash (#) that is ignored by R.
- Always set your working directory to the sub-folder containing your R script.
- Always run your code from an R script... **always!**

1.4.1 References

Brunsdon, Chris, and Lex Comber. 2018. *An Introduction to r for Spatial Analysis and Mapping (2e)*. Sage.

Comber, Lex, and Chris Brunsdon. 2021. *Geographical Data Science and Spatial Data Analysis: An Introduction in r*. Sage.

Harris, Richard. 2016. *Quantitative Geography: The Basics*. Sage.

Other good on-line *get started in R* guides include:

- The Owen guide (only up to page 28) : <https://cran.r-project.org/doc/contrib/Owen-TheRGuide.pdf>
- An Introduction to R - https://cran.r-project.org/doc/contrib/Lam-IntroductionToR_LHL.pdf
- R for beginners https://cran.r-project.org/doc/contrib/Paradis-rdebuts_en.pdf

1.5 Formative Tasks

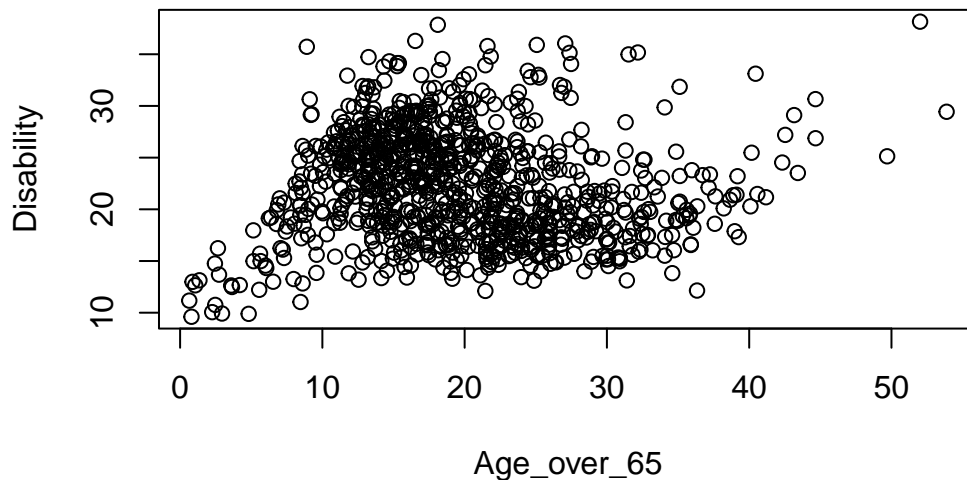
Task 1 From the district level dataset “merseyside.csv”, extract household information for the Liverpool and Wirral districts. The variables to be included are “Households”, “No_central_heating” and “Overcrowding”.

```
df <- read.csv("merseyside.csv")
df[c(2,5),c("District","Households","No_central_heating","Overcrowding")]
```

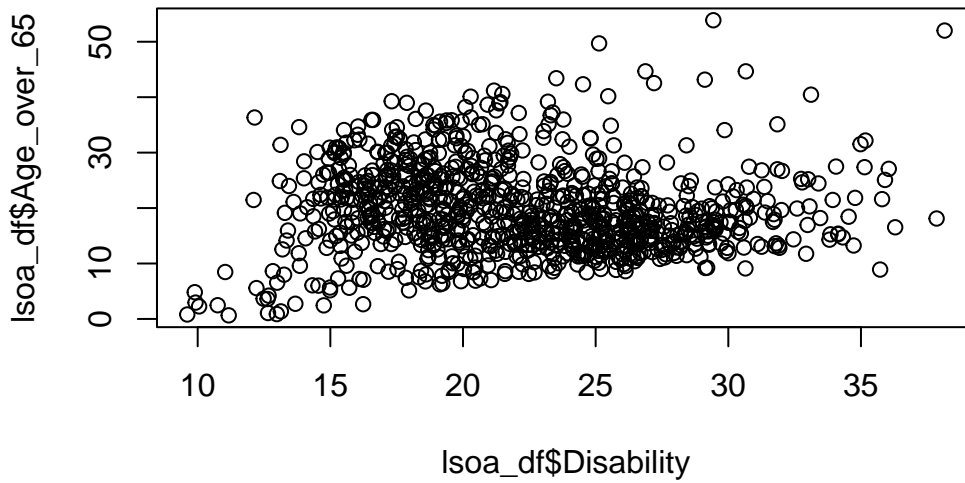
| | District | Households | No_central_heating | Overcrowding |
|---|-----------|------------|--------------------|--------------|
| 2 | Liverpool | 207491 | 4822 | 7352 |
| 5 | Wirral | 143253 | 2125 | 2355 |

Task 2 Use the dataset “merseyside_lsoa.csv”, plot the Disability against Age_over_65 from the data frame.

```
lsoa_df <- read.csv("merseyside_lsoa.csv")
plot(Disability~Age_over_65, data = lsoa_df)
```



```
# or
plot(lsoa_df$Disability, lsoa_df$Age_over_65)
```



Task 3 Use the district level dataset, how many households in total in Merseyside?

```
df <- read.csv("merseyside.csv")  
sum(df$Households)
```

```
[1] 620903
```

Task 4 Use the district level dataset, what is the overall proportion of the ageing population (age over 65) in Merseyside?

```
sum(df$Age_over_65)/sum(df$Population)
```

```
[1] 0.1920626
```

Task 5 Use the LSOA level dataset, what is the average proportion of the ageing population (age over 65) across all the neighbourhoods of Merseyside?

```
df$ageing_rate = df$Age_over_65 / df$Population * 100  
mean(df$ageing_rate)
```

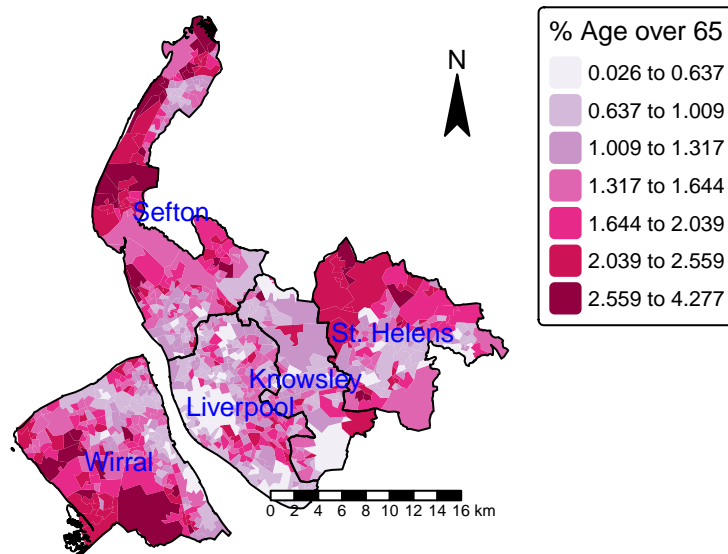
```
[1] 19.59824
```

Task 6 Create a map showing the spatial distribution of the proportion of ageing population (age over 65) over LSOAs in Merseyside? (use Jenks classification of 7 categories).

```
lsoa$ageing_rate = lsoa$Age_over_65 / lsoa$Population * 100

tm_shape(lsoa) +
  tm_fill(fill = "ageing_rate",
          fill.scale = tm_scale(values="PuRd",style = "jenks",n=7),
          fill.legend = tm_legend(title = "% Age over 65", title.size = 0.8)
  ) +
  tm_layout(main.title = "Merseyside",
            main.title.size=1.2,
            frame = FALSE) +
  tm_compass(position = c("right", "top")) +
  tm_scalebar(position = c("right", "bottom")) +
  tm_shape(merseyside) + # Add another spatial layer (Merseyside boundary)
  tm_borders(col = "black", lwd = 1) + # Draw the boundaries with black lines of width 1
  tm_text("LAD25NM",col = "blue",size = 0.8)
```

Merseyside



2 Lab: Exploratory Data Analysis - UK Election

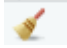
2.1 Overview

This week's practical session will draw upon the UK 2024 constituency election dataset. We will revisit some libraries and functions we have used last week, but also learn to do our first exploratory data analysis by fundamental R coding. We will do:

- Loading and examining tabular data (like spreadsheets) and geographical data into R (same as Week 1)
- Exploratory Data Analysis (or EDA) of numeric variables and categorical variables
- Using histogram, boxplot, barplot to understand the distribution of variables
- Variable interactions, particularly cross-tabulation and between-group comparison

You may wish to recap this week's lecture: [Lecture 02.pptx](#)

2.2 Clear the decks

- For this Week 2 session, create a sub-folder called **Week2** in your **ENVS162** folder on your M-Drive.
- Open RStudio
- Open a new R Script for your Week 2 work, rename it as **Week2.R** and save it in your newly created Week 2 folder, under M drive -> **ENVS162** folder. This is exactly the step we did in Week 1, and we will do this every week to Week 5.
- Check whether there is any previous left dataframes in your RStudio in the upper-right side Environment pane. You can always use the  to clear all the dataframes in your environment and make it all clean. For the same aim, you can run the below code:

```
rm(list = ls())
```

This command will clear RStudio's memory, removing any data objects that you have previously created.

2.3 Open libraries

In Week 1 we have installed essential R package `tidyverse`, `sf`, and `tmap`. Remember if any package has been installed, then we don't need to re-install them. Instead, we use `library()` command to import and use them.

As ever, when you start a new session in RStudio, you need to load the packages you wish to use into memory. Similarly, there are `tidyverse`, `tmap` and `sf` packages we've used last week.

```
library(tidyverse)
library(sf)
library(tmap)
```

If R returns Error: there is no package called '***'. Then it means that the package '***' has not been installed in the PC you current use. Therefore you need to install them first. Switch back to Week 1 instruction - Getting set up with RStudio - Your first R code - Package part to refresh yourself how to do this.

2.4 Parliamentary Constituency Data

2.4.1 Load the dataset

In 2024 the UK held a general election. Download the file ***uk_constituencies_2024.csv*** from our [Canvas module page Week 2](#). Save the dataset in your M drive - ENVS162 - Week 2 folder, alongside the Week 2. R script. Read in the dataset exactly in the same way as we did in Week 1 - insert the below code line and run it:

```
pc_data <- read.csv("uk_constituencies_2024.csv", stringsAsFactors = TRUE)
```

The datasheet captures a range of information relating to this election, and to the nature of each parliamentary constituency. By using `read.csv()` command, we use R to store the dataset in a dataframe called `pc_data` (short for Parliamentary Constituency data).

The `pc_data` dataset should appear in your RStudio Environment Pane on the upper right part, indicating that it has now been loaded into memory, and is available for analysis.

For your information only, the `pc_data` dataset has been assembled by combining information from the following sources.

- *HoC-GE2024-results-by-constituency.xlsx*

Source: Cracknell et al (2024) General election 2024 results, *Research Briefing*, House of Commons Library. <https://commonslibrary.parliament.uk/research-briefings/cbp-10009/>

- *Demographic-data-for-new-parliamentary-constituencies-May-2024.xlsx*

Source: House of Commons Demographic data for Constituencies <https://commonslibrary.parliament.uk/demographic-data-for-new-parliamentary-constituencies/>

- *NatCen Constituency Data_20 June 2024.xlsx*

Source: National Centre for Social Research (2024) Parliamentary constituency look-up, <https://natcen.ac.uk/constituency-look-up> (date accessed: 07-02-2025)

- *nomis_2025_01_14_101749.xlsx*

Source: Income data from Annual Survey of Hours and Earnings (ASHE), collected by the Office for National Statistics. <https://www.nomisweb.co.uk/datasets/asher>

2.4.2 Familiar with the dataset and variable types

In the `pc_data` dataset each row represents a different UK Parliamentary Constituency. Use the `View()` command to familiarise yourself with the variables contained in the dataset.

```
View(pc_data)
```

Use the `nrow()` or `dim()` command to find out how many Parliamentary Constituencies (and therefore MPs) there are in the UK.

```
nrow(pc_data)
```

```
dim(pc_data)
```

Whenever we try to understand the variables in one dataset, the first question to ask ourselves is: “are they *continuous* or *categorical*?”

Continuous variables are numeric measures of some quantity, such as a count or percentage or a precise value. E.g. number of valid votes; % of persons unemployed etc.

In contrast, categorical variables simply group observations into categories or ranges. E.g. name of the winning party; age group etc.

Explore the structure of the data table using the `str` function and have examined it by `head` functions.

```
str(pc_data)

head(pc_data)
```

We can see that we have numeric data in integers (`int`) form (these are counts or whole numbers) and continuous (`num`) form, and the character variables (text) have been converted to **factors**. For each of these data types we can generate numeric and visual summaries and we can also see how they interact with each other.

The `pc_data` dataset contains three basic sets of information about each Parliamentary Constituency:

- **constituency identifiers** - `gss_code` and `pc_name`
- **population information** for each constituency, ranging from the total population and number of households through to the % in various categories to information about local house prices, salaries and crime rates
- **2024 election results** ranging from the winning MP and party through to the size of the electorate and vote turnout, and the share of votes received by each party

2.4.3 Exploratory Data Analysis (EDA)

2.4.3.1 Numeric variables

You can use the `pc_data` dataset to extract some headline results from the 2024 General Election. Starting with the simplest case, the distribution of a single numeric variable whether continuous or count based can be examined numerically using the `summary` function.

```
#valid votes in constituency
summary(pc_data$valid_votes)
```

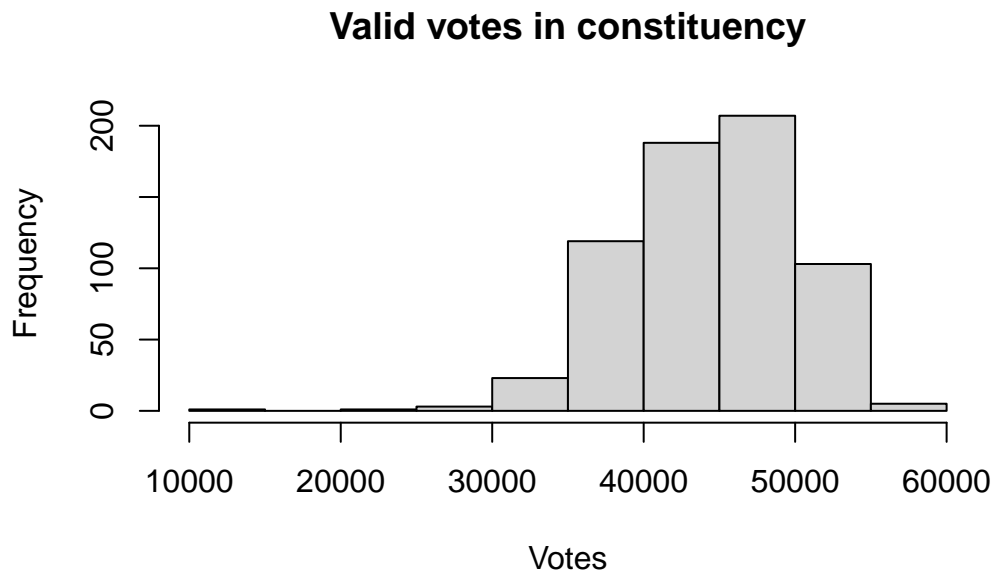
| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
|-------|---------|--------|-------|---------|-------|
| 13528 | 40397 | 44628 | 44322 | 48607 | 57744 |

```
# percentage of White British in constituency
summary(pc_data$pct_White_British)
```

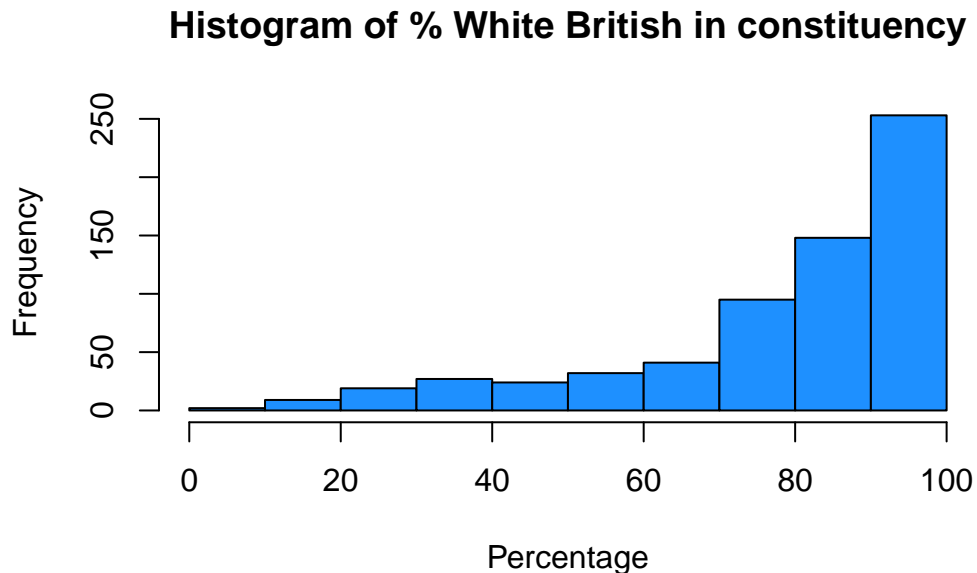
| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
|-------|---------|--------|--------|---------|--------|
| 8.937 | 71.534 | 86.208 | 78.006 | 92.843 | 98.664 |

A visual approach is more intuitive. The code below plots histograms of the two variables, using the `hist` function. The logic under-pinning a histogram is that a continuous variable (in this case `valid_votes` and `pct_White_British`) is temporarily regrouped into categories, using an equal-interval approach. The number of observations (in this case, constituencies) that fall into each equal-interval category then determine the height of each column in the histogram. The comments after each line shall inform you what the `main` and `xlab` in the function mean:

```
# histograms
hist(pc_data$valid_votes,
     main = "Valid votes in constituency", #change chart title
     xlab = "Votes")                      #change x axis label
```



```
hist(pc_data$pct_White_British,
     main = "Histogram of % White British in constituency",
     xlab = "Percentage",
     col = "dodgerblue") #change bar color to dodgerblue
```



You may noticed that `valid_votes` has a relatively normal, bell-shaped distribution whereas the `pct_White_British` variable is left skewed (negatively) distribution.

We can examine the how these distributions relate to central tendencies (mean, median) and spread, using standard deviations for means and the IQR (Inter-Quartile Range) for medians.

From the numeric summary above, the mean of `pc_data$valid_votes` is 44,322. We can determine the spread around this value by calculating the standard deviation for our sample , as is returned by the `sd` function in R:

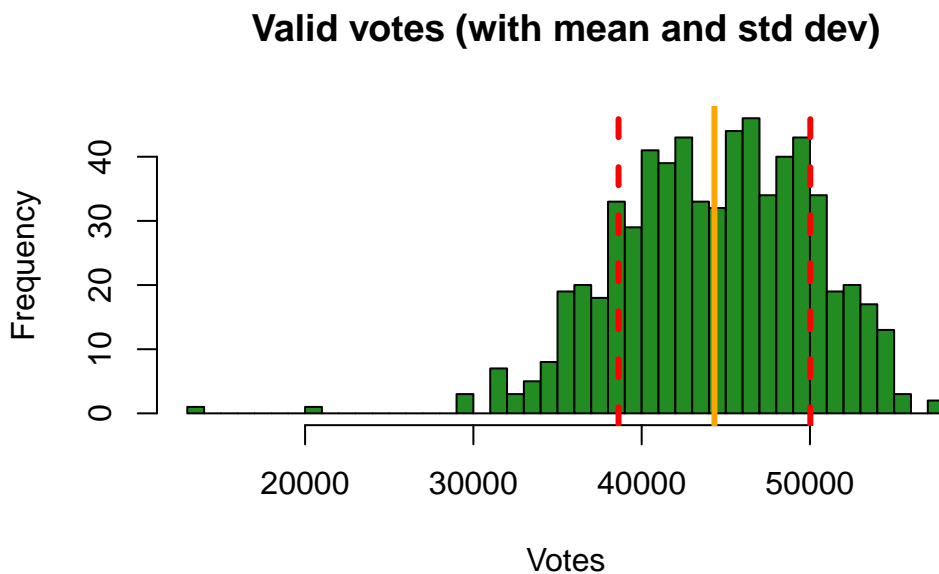
```
sd(pc_data$valid_votes)
```

```
[1] 5697.665
```

For a normal distribution, about 68% of observations lie within 1 standard deviation of the mean, and about 95% lie within 2 standard deviations of the mean. So this suggest that 68% of the valid votes are within 5,698 votes of the mean of 44,322 votes, i.e. 38,624 (44322-5698) and 50020 (44322+5698).

We can augment the histogram of the `valid_votes` variable with this information, requesting R to increase the intervals to 50, and using the `abline` function to add lines to create the figure below:

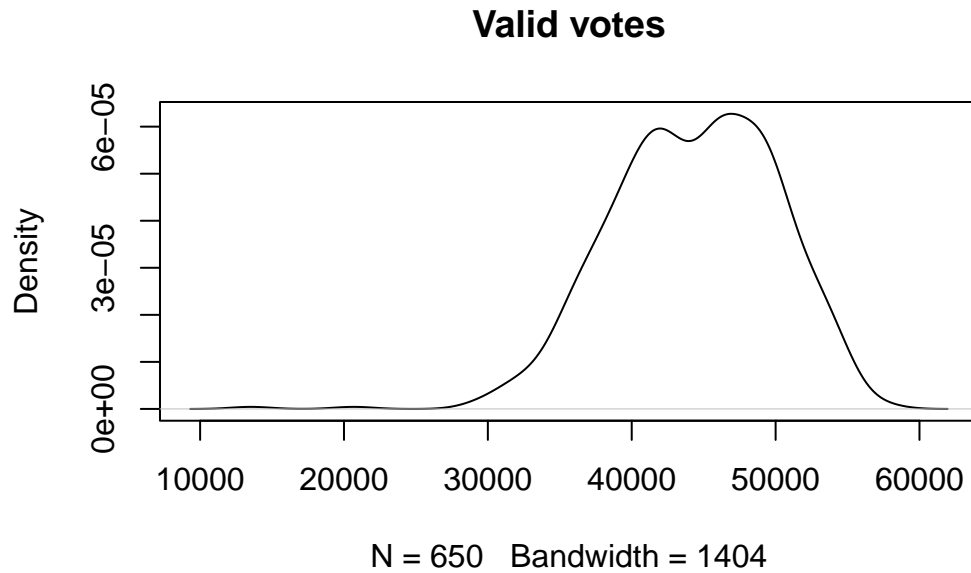
```
# histogram
hist(pc_data$valid_votes, col="forestgreen", main="Valid votes (with mean and std dev)",
     breaks = 50, xlab="Votes")
# calculate and add the mean
mean_val = mean(pc_data$valid_votes)
abline(v = mean_val, col = "orange", lwd = 3)
# calculate and add the standard deviation lines around the mean
sdev = sd(pc_data$valid_votes)
# minus 1 sd
abline(v = mean_val-sdev, col = "red", lwd = 3, lty = 2)
# plus 1 sd
abline(v = mean_val+sdev, col = "red", lwd = 3, lty = 2)
```



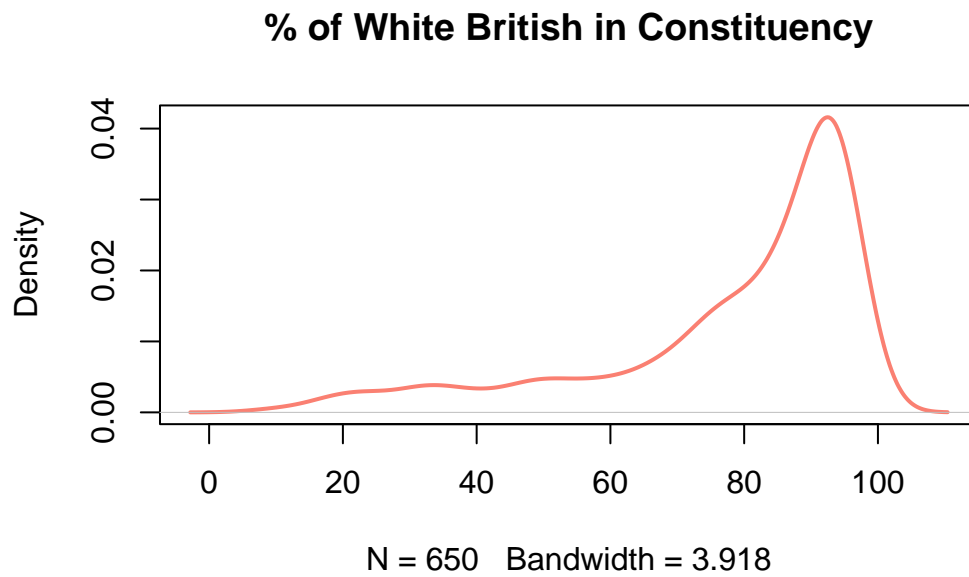
This histogram shows the variable `valid_votes` with the orange solid line as the mean, and dashed red lines as the standard deviation. Note that in the call to `abline` above, note the specification of different line types (`lty`) and line widths (`lwd`). Later on, you could explore these as described here.

We can also use density curve to present the variable distribution:

```
plot(density(pc_data$valid_votes),  
     main = "Valid votes")
```

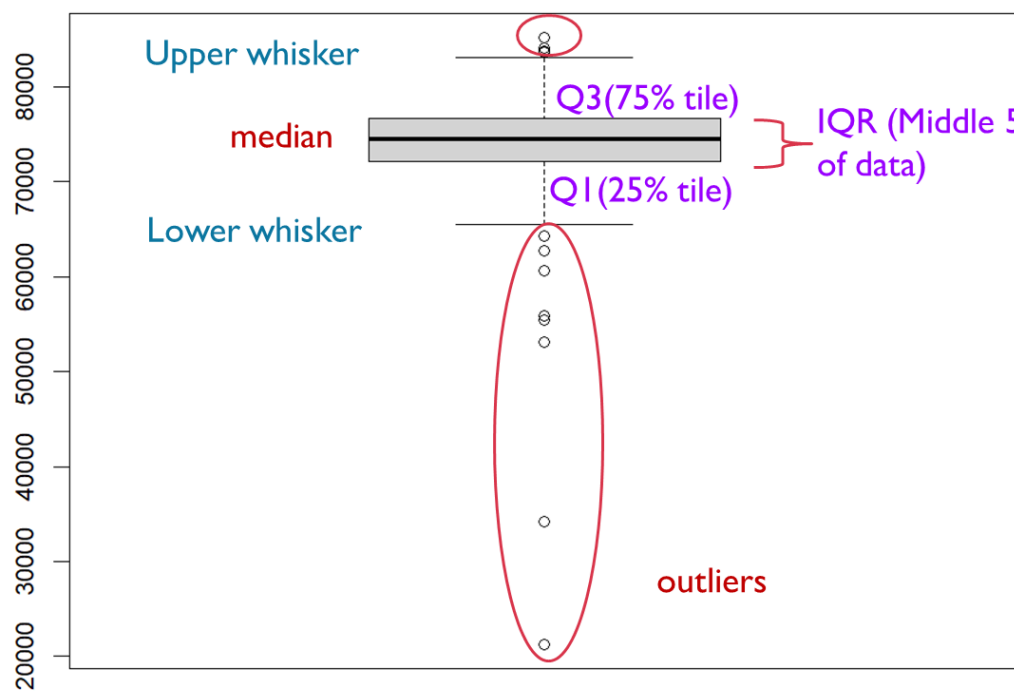


```
plot(density(pc_data$pct_White_British),  
     main = "% of White British in Constituency",  
     col="salmon",  
     lwd=2)
```



Boxplots show the same information but here we can see a bit more of the nature of the distribution.

Recap Week 2 lecture, the dark line shows the median, the box represents the interquartile range (from the 1st to the 3rd quartile), the whiskers extend to the most extreme non-outlying values, and the dots indicate outliers.



Roughly 99% of data

IQR = Interquartile Range

The box = IQR = middle 50% of data

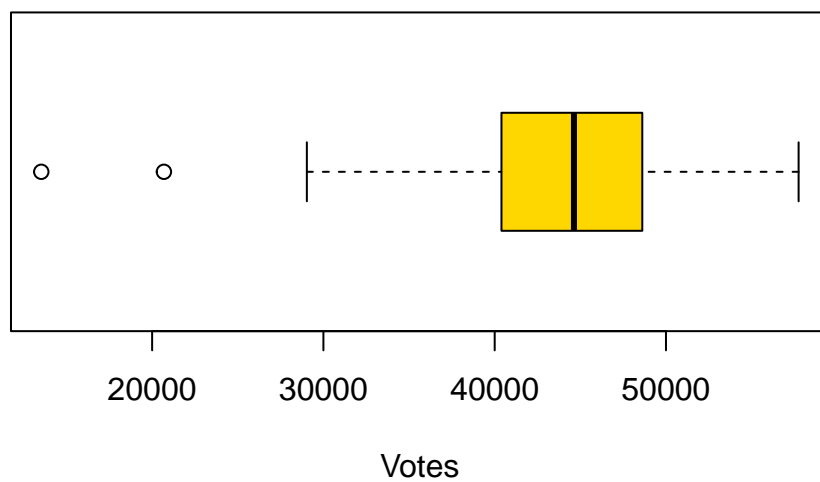
Lower whisker = smallest 25% of data

Upper whisker = largest 25% of data

Now let's check out the boxplots of these two variables:

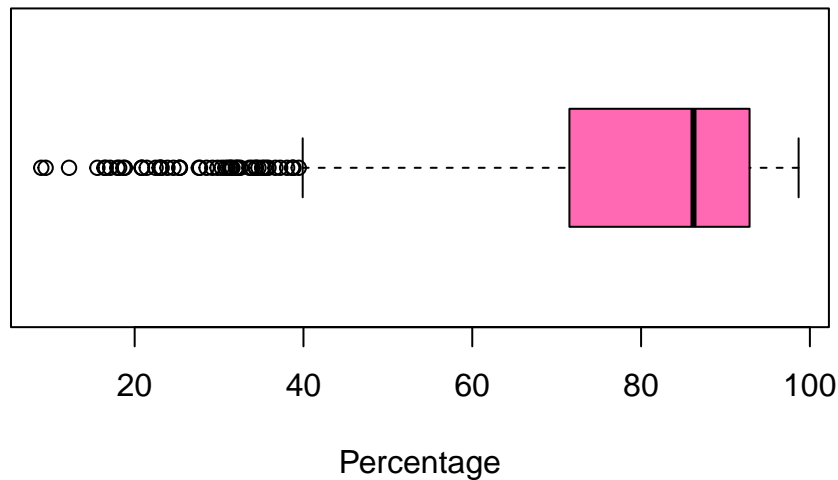
```
boxplot(pc_data$valid_votes,
        horizontal=TRUE,
        main = "Valid votes",
        xlab='Votes',
        col = "gold")
```


Valid votes



```
boxplot(pc_data$pct_White_British,  
        horizontal=TRUE,  
        main = "% of White British",  
        xlab='Percentage',  
        col = "hotpink")
```

% of White British



1. Is the median line in the centre of the box? If the median line is not in the middle of the box, it means the data are skewed, with greater spread on one side of the median.
 - Median closer to the bottom of the box -> right-skewed (positively skewed) -> some big values pulling up
 - Median closer to the top of the box -> left-skewed (negatively skewed) -> some small values dragging down
2. Are the whiskers the same length? This indicates skewness in the distribution. The data have a longer tail on the side of the longer whisker, meaning values are more spread out in that direction.
3. Are there any outliers? Who has more? Outliers clustered on one side indicate that extreme observations occur predominantly in one tail of the distribution. More outliers indicates the variable has more extreme or unusual values, possibly a heavier-tailed distribution. But this doesn't mean the variable is bad or more variable overall.
4. The size of the boxes? The box represents the interquartile range (IQR), which is the middle 50% of the data as we call them typical. A larger box indicates greater variability in the middle 50% of the variable; a smaller box suggests that values are more tightly clustered around the median.

Now, you should be able to compare how the two kinds of distribution are shown in the the boxplots with a trained eye: The `pct_White_British` variable is left-skewed, indicating that

more values are spread towards lower percentages of White British in the constituency. The box is also larger than that of `valid_votes`, which suggests greater variability in the central 50% of the data. In addition, the longer lower whisker and the presence of more outliers on the lower end further reinforce the left-skewed distribution.

In summary, numeric variable distributions, of counts and continuous data, should be investigated as an initial step in any data analysis. There are a number of metrics and graphical functions (tools) for doing this including `summary()`, `hist()` `plot(density())` and `boxplot()`.

2.4.3.2 Categorical variables

Some of the character variables could be considered as **categorical**, representing a grouping or classification of some kind, as described above. In these cases we are interested in the count or frequency of each class in the classification, which we can examine numerically or graphically.

The simplest way to examine classes is to put them into a table of counts. The `table` function is very useful and in the code below it is applied to one of the categorical variables in the survey data:

So firstly we use the `table()` command to find the number of MPs elected to each Party. [Hint: use the `first_party` variable].

```
table(pc_data$first_party)
```

| | | | | |
|-----------|--------------|-------------------------|----------------|-------------|
| Alliance | Conservative | DUP | Green | Independent |
| 1 | 121 | 5 | 4 | 6 |
| Labour | Lib Dems | Plaid Cymru | Reform | SDLP |
| 411 | 72 | 4 | 5 | 2 |
| Sinn Fein | SNP | Speaker Ulster Unionist | Unionist Voice | |
| 7 | 9 | 1 | 1 | 1 |

These can be made a bit more tabular in format with the `data.frame` function, which takes the `table` operation as its input:

```
data.frame(table(pc_data$first_party))
```

| | Var1 | Freq |
|---|--------------|------|
| 1 | Alliance | 1 |
| 2 | Conservative | 121 |
| 3 | DUP | 5 |

| | | |
|----|-----------------|-----|
| 4 | Green | 4 |
| 5 | Independent | 6 |
| 6 | Labour | 411 |
| 7 | Lib Dems | 72 |
| 8 | Plaid Cymru | 4 |
| 9 | Reform | 5 |
| 10 | SDLP | 2 |
| 11 | Sinn Fein | 7 |
| 12 | SNP | 9 |
| 13 | Speaker | 1 |
| 14 | Ulster Unionist | 1 |
| 15 | Unionist Voice | 1 |

However, if we not only care about the count of MPs but also the proportion? Then we can make a good use of our tidyverse library to run the following code line.

```
pc_data %>%
  count(first_party) %>%
  mutate(pct = round(n / sum(n) * 100,1))
```

| | first_party | n | pct |
|----|-----------------|-----|------|
| 1 | Alliance | 1 | 0.2 |
| 2 | Conservative | 121 | 18.6 |
| 3 | DUP | 5 | 0.8 |
| 4 | Green | 4 | 0.6 |
| 5 | Independent | 6 | 0.9 |
| 6 | Labour | 411 | 63.2 |
| 7 | Lib Dems | 72 | 11.1 |
| 8 | Plaid Cymru | 4 | 0.6 |
| 9 | Reform | 5 | 0.8 |
| 10 | SDLP | 2 | 0.3 |
| 11 | Sinn Fein | 7 | 1.1 |
| 12 | SNP | 9 | 1.4 |
| 13 | Speaker | 1 | 0.2 |
| 14 | Ulster Unionist | 1 | 0.2 |
| 15 | Unionist Voice | 1 | 0.2 |

Here you are using two very useful functions in the library `tidyverse`.

First, the `count()` calculate the frequency of different categories in the `first_party` and use a new column `n` to store the frequencies - it actually do the same thing as above code, but better in presenting as a table;

Second, the `mutate()` function to assist use create a new column `pct` and fill in the value by the calculation `pct = n / sum(n) * 100`.

The `%>%` is used to link these two commands: `count()` and `mutate()`. You can insert `%>%` in your R script by using `ctrl + shift + M` for Windows and `Cmd + Shift + M` on Mac.

Therefore, when you run the code, you will see a table showing as three column: first part name, a new column automatically named as `n` by R after the `count()` function, and count of the party MPs, and a new column called `pct` and with values calculated by `n / sum(n) * 100`. We use `round()` function to keep only 1 digits for the `pct`.

We can also improve the code to make a better table presentation. You may find the comment text after each code line would be useful to understand what R has done to the `pc_data`.

```
#Calculate the frequency and percentage of different categories for "first_party"
pc_data %>%
  count(first_party) %>%
  mutate(pct = round(n / sum(n) * 100,1)) %>%
  arrange(desc(n)) %>% #sort the table by number of MPs from more to less
  setNames(c("First Party", "Number of MPs", "% of MPs")) #rename table column names
```

| | First Party | Number of MPs | % of MPs |
|----|-----------------|---------------|----------|
| 1 | Labour | 411 | 63.2 |
| 2 | Conservative | 121 | 18.6 |
| 3 | Lib Dems | 72 | 11.1 |
| 4 | SNP | 9 | 1.4 |
| 5 | Sinn Fein | 7 | 1.1 |
| 6 | Independent | 6 | 0.9 |
| 7 | DUP | 5 | 0.8 |
| 8 | Reform | 5 | 0.8 |
| 9 | Green | 4 | 0.6 |
| 10 | Plaid Cymru | 4 | 0.6 |
| 11 | SDLP | 2 | 0.3 |
| 12 | Alliance | 1 | 0.2 |
| 13 | Speaker | 1 | 0.2 |
| 14 | Ulster Unionist | 1 | 0.2 |
| 15 | Unionist Voice | 1 | 0.2 |

In this code chunk, we requested four command to the dataframe `pc_data`, and linked them with `%>%` :

1. `count()` function to summarises the data by counting the number of observations in each group;

2. `mutate()` function to create a new column named “pct” as we did above;
3. `arrange()` function sorts the rows, `desc(n)` function decending the norder of `n`, together they order the category with the highest counts first;
4. `setNames()` function renames the results.

Similarly, we can use the categorical variable `crime_rate` in the `pc_data` dataset to understand the crime status of all the constituencies:

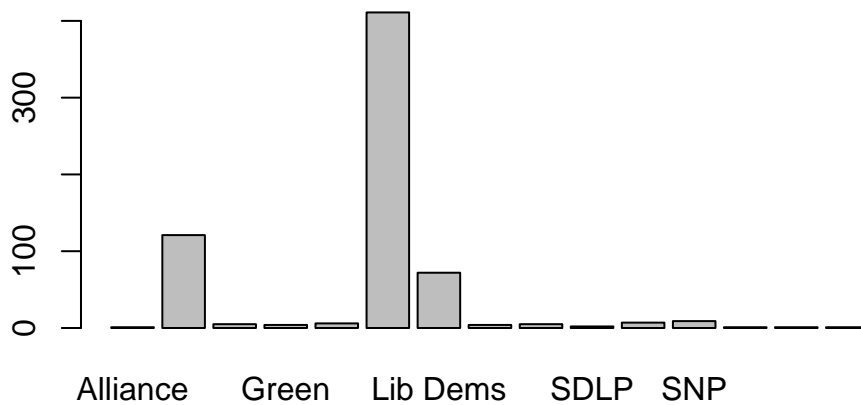
```
#Calculate the frequency and percentage of different categories for "crime_rate"
pc_data %>%
  count(crime_rate) %>%
  mutate(pct = round(n / sum(n)*100, 1)) %>%
  arrange(desc(n)) %>%
  setNames(c("Crime rate", "Number of Constituency", "% of Constituency"))
```

| | Crime rate | Number of Constituency | % of Constituency |
|---|-------------|------------------------|-------------------|
| 1 | Much higher | 118 | 18.2 |
| 2 | Higher | 115 | 17.7 |
| 3 | Average | 114 | 17.5 |
| 4 | Lower | 114 | 17.5 |
| 5 | Much lower | 114 | 17.5 |
| 6 | <NA> | 75 | 11.5 |

What you have learnt from the result table? It seems that there is a quite equally distribution across the five categories of crime rate, although there are 11.5% of the constituency are missing their crime rates.

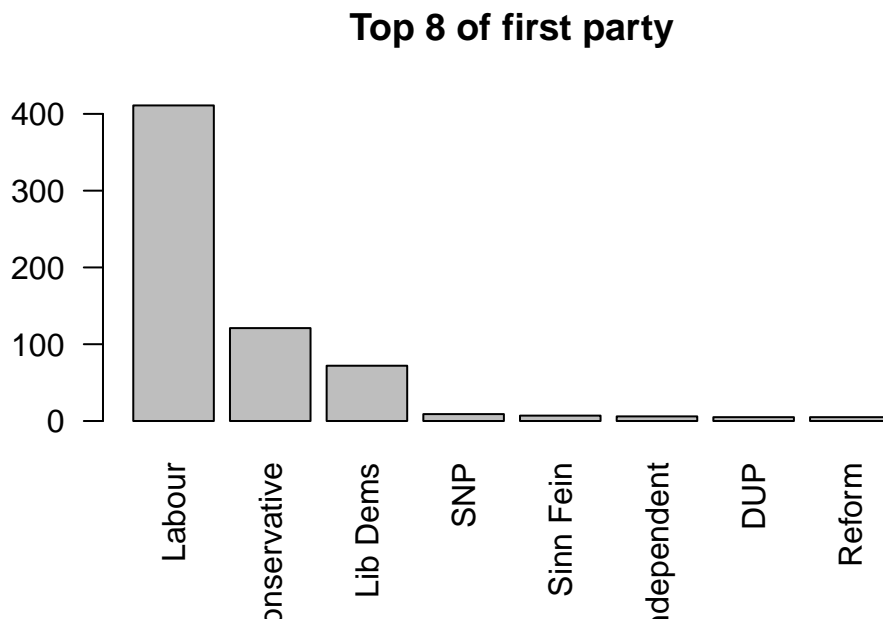
Categorical data can be visualised using bar plots of the tabularised data. The code below does this by creating a table, changing the names of the table and passing that to the `barplot` function:

```
#calculate the frequency of first party and save the result in table 'tab', present tab as a table
tab = table(pc_data$first_party)
barplot(tab)
```



It is very simple to get the barplot from the result of `table()`. But it may need some improvement. As you may noticed, there are many bars have rarely no values, and too crowd x axis labels makes some bar label can't able to show. Therefore, we probably don't need to show all the parties, but only the top 8 in terms of their winning constituencies.

```
#sorted tab by the frequency from highest to lowest, present the top 8 of the sorted tab
tab_sorted <- sort(tab,decreasing = TRUE)
barplot(head(tab_sorted,8),las = 2, main = "Top 8 of first party")
```



In this code chunk, we first use `sort()` to reorder the `tab` ranking the parties based on the counts from highest to lowest. You may ask R to show how `tab_sorted` looks like:

```
#show tab_sorted
tab_sorted
```

| | | | | |
|-------------|--------------|-------------------------|----------------|-------------|
| Labour | Conservative | Lib Dems | SNP | Sinn Fein |
| 411 | 121 | 72 | 9 | 7 |
| Independent | DUP | Reform | Green | Plaid Cymru |
| 6 | 5 | 5 | 4 | 4 |
| SDLP | Alliance | Speaker Ulster Unionist | Unionist Voice | |
| 2 | 1 | 1 | 1 | 1 |

Then, using `barplot()`, we plot the top 8 rows of `tab_sorted`, with `las = 2` used to rotate the axis labels so they are displayed vertically and remain readable.

2.4.3.3 Categorical to categorical: cross-tabulation

In EDA, it is also important to understand the relationship between variables. Now it is the time to do the variable interactions. Let's first start with interact one categorical variable to another categorical variable. In many situations, it can be called cross-tabulation.

The relationship between two sets of classes or categories can be explored using correspondence tables created by the `table` function. Here we can cross tabulate the two categorical variables that we have already familiar with: `first_party` and `crime_rate`:

If we want to examine how the distribution of crime-rate categories varies for each first party. The question can be: for each political party, how are its wins distributed across different crime-rate levels?

```
# cross-tabulation first_party vs crime_rate

table(pc_data$first_party, pc_data$crime_rate)
```

| | Average | Higher | Lower | Much higher | Much lower |
|-----------------|---------|--------|-------|-------------|------------|
| Alliance | 0 | 0 | 0 | 0 | 0 |
| Conservative | 25 | 6 | 44 | 0 | 41 |
| DUP | 0 | 0 | 0 | 0 | 0 |
| Green | 1 | 0 | 0 | 1 | 2 |
| Independent | 0 | 3 | 0 | 2 | 0 |
| Labour | 79 | 101 | 47 | 115 | 32 |
| Lib Dems | 6 | 2 | 20 | 0 | 38 |
| Plaid Cymru | 0 | 0 | 3 | 0 | 1 |
| Reform | 2 | 3 | 0 | 0 | 0 |
| SDLP | 0 | 0 | 0 | 0 | 0 |
| Sinn Fein | 0 | 0 | 0 | 0 | 0 |
| SNP | 0 | 0 | 0 | 0 | 0 |
| Speaker | 1 | 0 | 0 | 0 | 0 |
| Ulster Unionist | 0 | 0 | 0 | 0 | 0 |
| Unionist Voice | 0 | 0 | 0 | 0 | 0 |

Conversely, we can also examine how the distribution of first parties varies across different crime-rate categories. The question this time is: for each crime-rate category, how are different parties distributed?

```
# cross-tabulation crime_rate vs first_party

table(pc_data$crime_rate, pc_data$first_party)
```

| | Alliance | Conservative | DUP | Green | Independent | Labour | Lib Dems |
|---------|----------|--------------|-----|-------|-------------|--------|----------|
| Average | 0 | 25 | 0 | 1 | 0 | 79 | 6 |
| Higher | 0 | 6 | 0 | 0 | 3 | 101 | 2 |

| | | | | | | | |
|-------------|---|----|---|---|---|-----|----|
| Lower | 0 | 44 | 0 | 0 | 0 | 47 | 20 |
| Much higher | 0 | 0 | 0 | 1 | 2 | 115 | 0 |
| Much lower | 0 | 41 | 0 | 2 | 0 | 32 | 38 |

| | Plaid Cymru | Reform | SDLP | Sinn Fein | SNP | Speaker | Ulster | Unionist |
|-------------|-------------|--------|------|-----------|-----|---------|--------|----------|
| Average | 0 | 2 | 0 | | 0 | 0 | 1 | 0 |
| Higher | 0 | 3 | 0 | | 0 | 0 | 0 | 0 |
| Lower | 3 | 0 | 0 | | 0 | 0 | 0 | 0 |
| Much higher | 0 | 0 | 0 | | 0 | 0 | 0 | 0 |
| Much lower | 1 | 0 | 0 | | 0 | 0 | 0 | 0 |

| | Unionist Voice |
|-------------|----------------|
| Average | 0 |
| Higher | 0 |
| Lower | 0 |
| Much higher | 0 |
| Much lower | 0 |

What insights now you can learn from these cross-tabulation? We will examine methods for determining whether the cross-tabulated counts and their differences are significant (i.e. would not be expected by chance) in later weeks.

2.4.3.4 Continuous to categorical: compare between groups

We may also be interested in the exploring differences and similarities in the continuous variables associated with for each categorical class. This can be done using multiple boxplots.

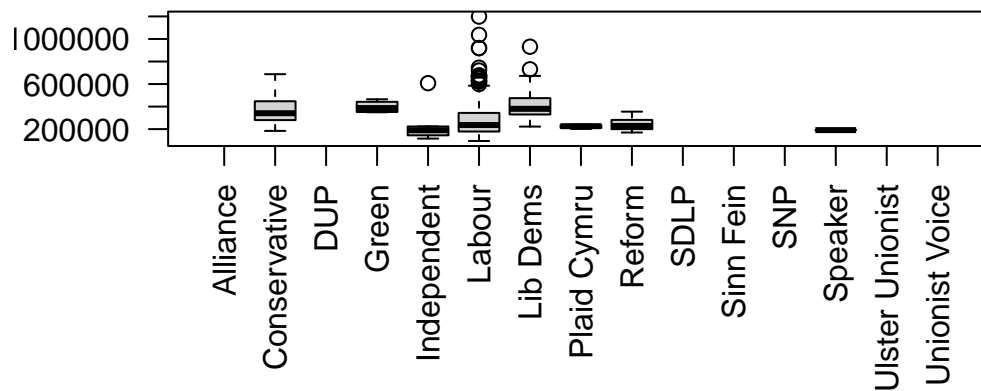
If we want to make boxplots of each party by the value median house price of the constituencies who vote for this party as their first party:

```
#create boxplots by median_house_price for each first_party

par(mar = c(10, 4, 4, 2)) # increase the margin of the chart for each side: bottom, left, top, right

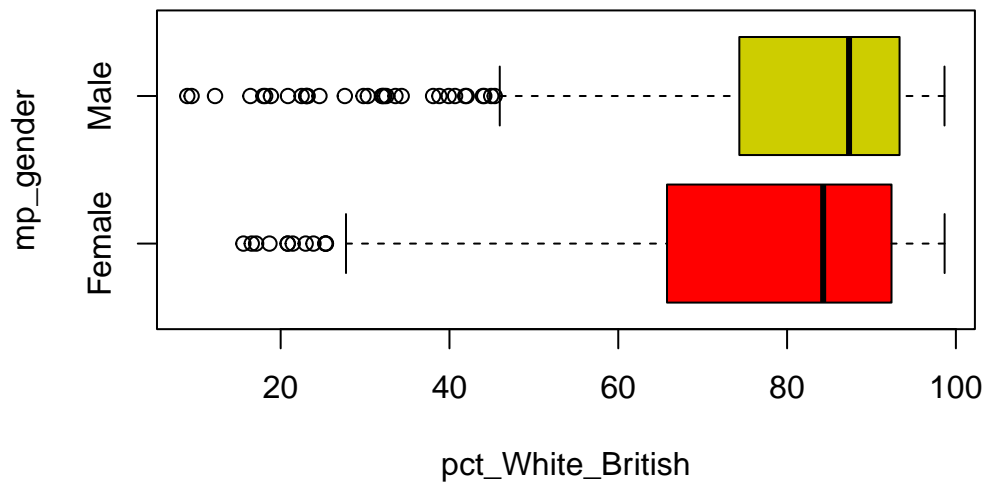
boxplot(median_house_price ~ first_party,
        data = pc_data,
        las = 2, #vertical present item label
        xlab="", #no x axis label
        ylab="", #no y axis label,
        main="Compare constituency median house price vs. first party")
```

Compare constituency median house price vs. first party



If we want to explore the % of White British in the constituencies and compare the distribution between different MP genders. The code below shows that, compared with constituencies represented by male MPs, those with female MPs generally have a lower proportion of White British population. This distribution also exhibits clear skewness towards lower percentages of White British residents.

```
boxplot(pct_White_British ~ mp_gender,
        data = pc_data, col=c("red", "yellow3"), #specify bar colors
        horizontal = TRUE) #horizontal the boxplot
```



We can do this numerically as well, but it is a bit more convoluted using the `with` and `aggregate` functions:

```
with(pc_data, aggregate(pct_White_British, by=list(mp_gender) , FUN=summary))
```

| | Group.1 | x.Min. | x.1st Qu. | x.Median | x.Mean | x.3rd Qu. | x.Max. |
|---|---------|-----------|-----------|-----------|-----------|-----------|-----------|
| 1 | Female | 15.587153 | 65.775762 | 84.276710 | 75.207471 | 92.374361 | 98.663686 |
| 2 | Male | 8.936846 | 74.343213 | 87.344901 | 79.907108 | 93.326348 | 98.652616 |

2.5 Make your own map for the election result

Having established how many MPs and votes each party got, it is time to look at the geography of the election outcome. To do this we need to link a set of digital boundaries for Parliamentary Constituencies with our `pc_data` dataset.

2.5.1 Read in Parliamentary Constituency Boundaries

Digital boundaries for the Parliamentary Constituencies used in the 2024 General Election can be found in the files *uk_constituencies_2024.gpkg* from the Cavans module page. Download and save it in the Week 2 folder as well.

As Week 1, we use `st_read()` function from `library(sf)` to read in this geographical boundary dataset.

```
#read the boundaries as a spatial dataset  
  
pc_map <- st_read("uk_constituencies_2024.gpkg")
```

```
Reading layer `uk_constituencies_2024' from data source  
  `C:\Users\ziye\Documents\GitHub\quant\labs\uk_constituencies_2024.gpkg'  
  using driver `GPKG'  
Simple feature collection with 650 features and 9 fields  
Geometry type: MULTIPOLYGON  
Dimension:      XY  
Bounding box:   xmin: 191.9359 ymin: 7423.9 xmax: 655599.6 ymax: 1218591  
Projected CRS: OSGB36 / British National Grid
```

2.5.2 Inspect the spatial dataset

Use the `names()` or `str()` to know the contents, or as above use `View()` to open and check. Check by yourself the number of rows and columns of the map data:

```
str(pc_map)
```

```
View(pc_map)
```

The `pc_map` dataset contains the standard set of Parliamentary Constituency boundaries:

```
#make a map of constituency  
tm_shape(pc_map) + #map a spatial data  
  tm_polygons() #map it as polygons
```

or we can make a colorful map by using different color for different regions:

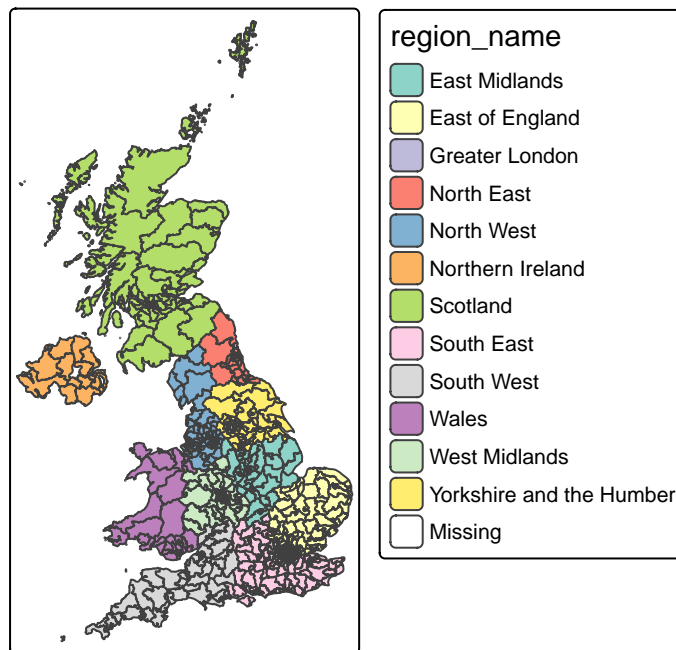
```
#make a map of constituency and color each polygon based on "region_name", using different c  
  
tm_shape(pc_map) + #map a spatial data  
  tm_polygons("region_name",  
    palette="Set3")
```

-- tmap v3 code detected -----

[v3->v4] ``tm_tm_polygons()``: migrate the argument(s) related to the scale of the visual variable ``fill`` namely 'palette' (rename to 'values') to `fill.scale = tm_scale(<HERE>)`.

[cols4all] color palettes: use palettes from the R package cols4all. Run ``cols4all::c4a_gui()`` to explore them. The old palette name "Set3" is named "brewer.set3"

Multiple palettes called "set3" found: "brewer.set3", "hcl.set3". The first one, "brewer.set3"



2.5.3 Link boundaries to pc_data

In order to map the election results contained in the `pc_data` dataset, we need to join it to a set of digital boundaries using the `left_join()` command - you should have already familiar with this from Week 1.

In your inspection of the `pc_data` and `pc_map` datasets, you may have noticed that they all have two variables in common. The first is a unique identifier for each Parliamentary Constituency: `gss_code`. The second is the name of the constituency: `pc_name`.

We can use these two variables to first link the `pc_data` dataset to the standard map:

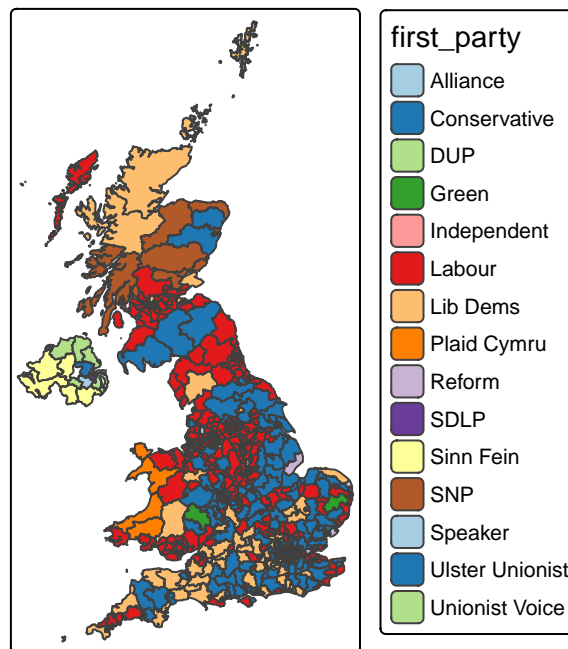
```
#left join pc_data to pc_map, joining when gss_code from pc_map equals to pc_name in pc_data
pc_map_new <- left_join(pc_map, pc_data, by = c("gss_code", "pc_name"))
```

As ever, having created a new dataset, use the `dim()`, `str()`, `names()` and `View()` commands to check its contents are as you would expect.

2.5.4 Map the election result


Having joined the `pc_data` dataset with a set of digital boundaries, it becomes a simple matter to map the election results using the mapping skills covered in Week 1:

```
#make a map for the pc_map_new, fill the colors for each polygon based on "first_party", using
tm_shape(pc_map_new) + #map a spatial data
  tm_polygons(fill = "first_party",
              fill.scale = tm_scale(values="Paired")) #map it as polygons, use different colors
```



This time, let's change the mode of `tmap` by using `tmap_mode()` from default "plot" to "view" for an interactive map:

```
# make the map interactive
tmap_mode("view")
tm_shape(pc_map_new) +
  tm_polygons("first_party")
```

This time, in your right-bottom pane, the map should be plotted in Viewer tab as an interactive map. You can zoom in/out to explore your map, for a better view, you can click the  to open a webpage.

You may need to switch `tmap_mode()` back to “plot” for a static map making:

```
#switch back to plot static mode for later use

tmap_mode("plot")
```

i tmap mode set to "plot".

Now what pattern you can observe from the interactive map you just made?

2.6 Formative tasks

Task 1 Write code to get how many columns and rows in the UK constituency boundary dataset `pc_map`?

```
pc_map <- st_read("uk_constituencies_2024.gpkg")
```

```
Reading layer `uk_constituencies_2024' from data source
  `C:\Users\ziye\Documents\GitHub\quant\labs\uk_constituencies_2024.gpkg'
  using driver `GPKG'
Simple feature collection with 650 features and 9 fields
Geometry type: MULTIPOLYGON
Dimension:      XY
Bounding box:   xmin: 191.9359 ymin: 7423.9 xmax: 655599.6 ymax: 1218591
Projected CRS:  OSGB36 / British National Grid
```

```
ncol(pc_map)
```

```
[1] 10
```



```
nrow(pc_map)
```

```
[1] 650
```

```
#both col and row  
dim(pc_map)
```

```
[1] 650 10
```

Task 2 Using the UK constituency boundary dataset `pc_map`, write code to get descriptive summary the area (variable: `sq_km`) of all the constituencies.

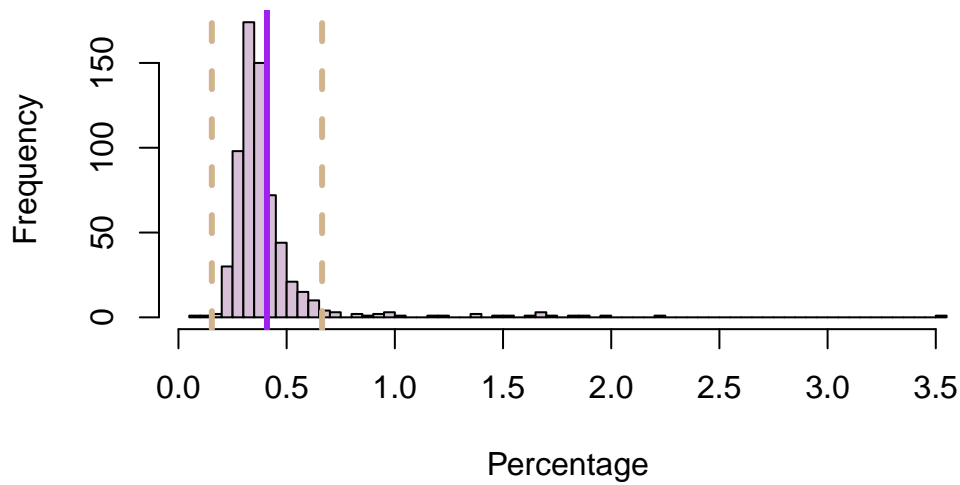
```
summary(pc_map$sq_km)
```

| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
|------|---------|--------|--------|---------|----------|
| 6.80 | 33.65 | 106.45 | 375.01 | 351.98 | 11634.40 |

Task 3 Write some codes to plot a histogram of the `pct_invalid_votes` variable in the constituency election dataset, with lines showing the mean and the standard deviation around the mean. Try add breaks into the function and change breaks from 10, 20 to 50 and see how the histogram changed.

```
pc_data <- read.csv("uk_constituencies_2024.csv", stringsAsFactors = TRUE)  
  
# histogram  
hist(pc_data$pct_invalid_votes, col = "thistle", main = "Invalid votes in constituency", xlab = "pct_invalid_votes")  
# calculate and add the mean  
mean_val = mean(pc_data$pct_invalid_votes, na.rm = T)  
abline(v = mean_val, col = "purple", lwd = 3)  
# calculate and add the standard deviation lines around the mean  
sdev = sd(pc_data$pct_invalid_votes, na.rm = T)  
abline(v = mean_val-sdev, col = "tan", lwd = 3, lty = 2)  
abline(v = mean_val+sdev, col = "tan", lwd = 3, lty = 2)
```

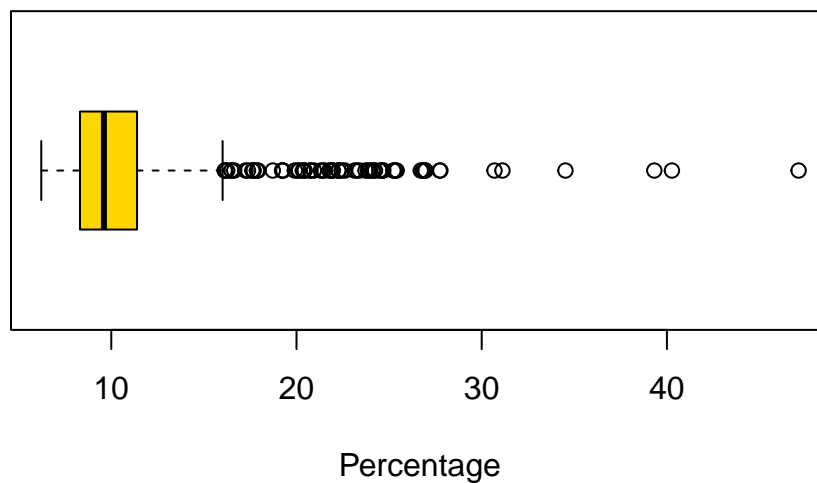
Invalid votes in constituency



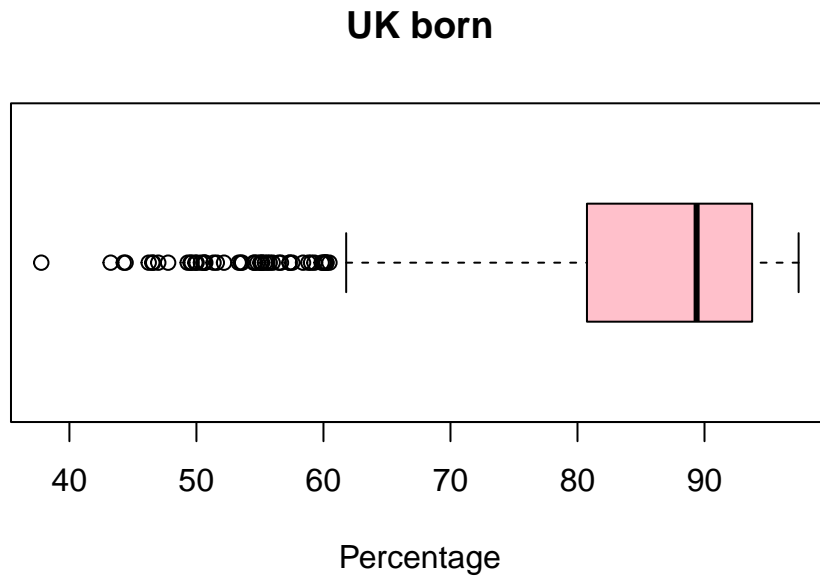
Task 4 Write codes to create boxplots for `pct_in_migration`, `pct_UK_born` and `pct_owned` (owning upright household):

```
boxplot(pc_data$pct_in_migration, horizontal=TRUE, main = "In migration", xlab='Percentage', c
```

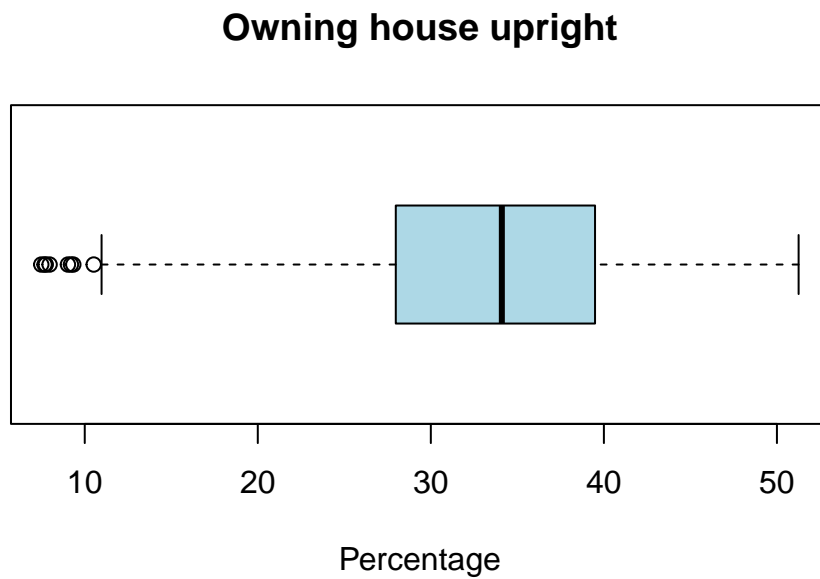
In migration



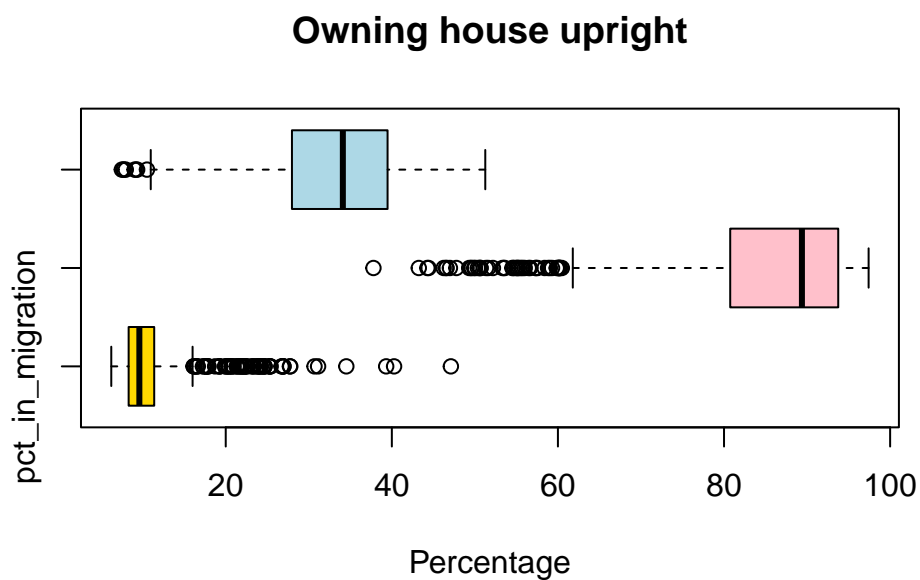
```
boxplot(pc_data$pct_UK_born ,horizontal=TRUE, main = "UK born", xlab='Percentage', col = "pink",
```



```
boxplot(pc_data$pct_owned,horizontal=TRUE, main = "Owning house upright", xlab='Percentage',
```



```
#to compare them together
boxplot(pc_data[,c("pct_in_migration",
                  "pct_UK_born",
                  "pct_owned"
                )],
        horizontal=TRUE,
        main = "Owning house upright",
        xlab='Percentage',
        col = c("gold","pink","lightblue")
      )
```



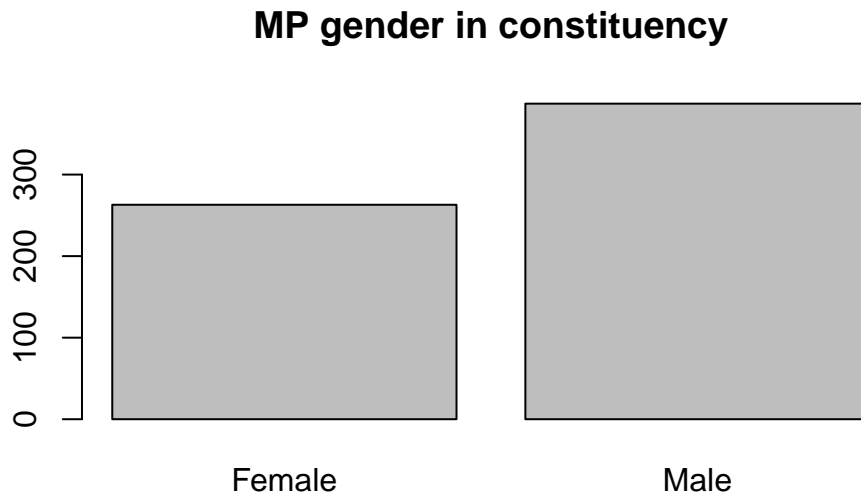
Task 5 Write codes to summary the counts of MPs in each gender (this is in the `mp_gender` variable), presenting the table with percentage and showing the barplot.

```
pc_data %>%
  count(mp_gender) %>%
  mutate(pct = round(n/sum(n)*100,1))
```

| | mp_gender | n | pct |
|---|-----------|-----|------|
| 1 | Female | 263 | 40.5 |
| 2 | Male | 387 | 59.5 |

```
tab = table(pc_data$mp_gender)

barplot(tab,main = "MP gender in constituency")
```



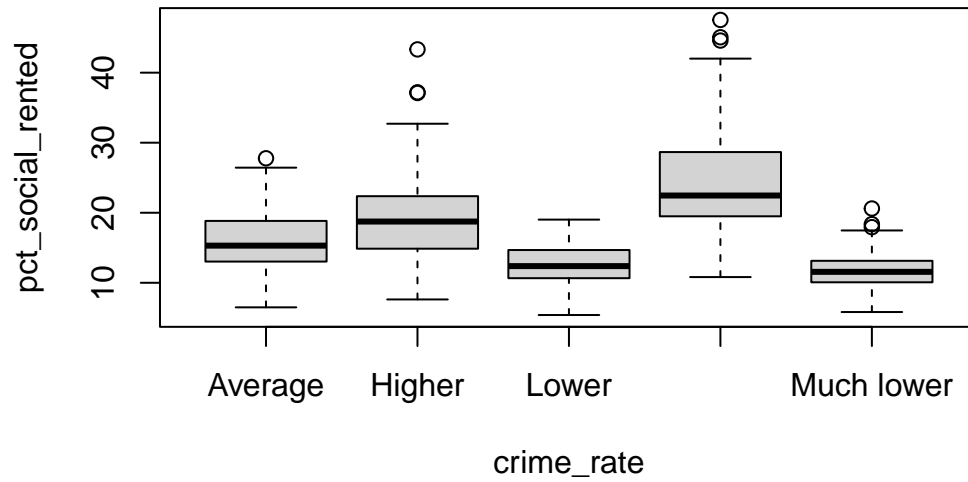
Task 6 Cross tabulation `region_name` to `mp_gender` by using the newly created `pc_map_new` dataset, which joined by the constituency boundary and election dataset.

```
table(pc_map_new$region_name, pc_map_new$mp_gender)
```

| | Female | Male |
|--------------------------|--------|------|
| East Midlands | 20 | 27 |
| East of England | 16 | 45 |
| Greater London | 38 | 37 |
| North East | 12 | 15 |
| North West | 31 | 42 |
| Northern Ireland | 5 | 13 |
| Scotland | 20 | 37 |
| South East | 39 | 52 |
| South West | 19 | 39 |
| Wales | 15 | 17 |
| West Midlands | 25 | 32 |
| Yorkshire and the Humber | 23 | 31 |

Task 7 Compare between `crime_rate` with the `pct_social_rented` in constituency election dataset. What pattern you can learn from your boxplot?

```
boxplot(pct_social_rented ~ crime_rate , data = pc_data)
```



3 Lab: Introductory Statistics - Happiness around the world

3.1 Overview

We have used Census data for the past two weeks. Census data aim to collect detailed information about every person and household in the UK, particularly on personal characteristics, household and housing, work and education, health, migration, and so on. The aim of such collection is to support plan public services and infrastructures.

Different from census (which counts the entire population), in many cases, we need to use surveys to obtain timely, cost-effective and probably more specific in-depth information from people. Therefore, this week's practical session will draw upon two survey datasets:

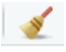
- UK Family Resource Survey
- World Value Survey

We will use these survey datasets to investigate the sample more thoroughly and employ statistical approaches to determine its representativeness of the population. We will use both numerical and categorical variables to calculate the sample mean, Standard Error and Confidence Intervals in different sample size.

You may wish to recap this week's lecture: [Lecture 03](#)

3.2 Prepare your working environment

- For this Week 3 session, create a sub-folder called Week3 in your ENVS162 folder on your M-Drive. This is exactly the step we did in Week 1 and 2 and we will do this every week to Week 5.
- Download this week's practical datasets from Canvas Week 3: *family_resource_survey.csv*, *world_value_survey.csv* and *world_map.geojson*. Save all these three datasets in the Week 3 folder you just created.
- Open RStudio

- Open a new R Script for your Week 3 work, rename it as Week3.R and save it in your newly created Week 3 folder, under M drive -> ENVS162 folder.
- Check whether there is any previous left dataframes in your RStudio in the upper-right side Environment pane. You can always use the to clear all the dataframes in your environment and make it all clean. For the same aim, you can click the icon , or you can run the below code:

```
rm(list = ls())
```

This command and also the brush icon can both clear RStudio's memory, removing any data objects that you have previously created.

3.3 Load libraries and familiar with survey data

Exactly as what we have done for Week 1 and 2, before we start to do anything in R, we first need to load the essential libraries as all the functions/commands are packed in these libraries. For this week, we will still rely on `tidyverse`, `tmap` and `sf`.

```
library(tidyverse)
library(sf)
library(tmap)
```

Recall that the data can be loaded into RStudio using the `read.csv` function:

```
# use read.csv to load a CSV file
df <- read.csv("family_resource_survey.csv", stringsAsFactors = TRUE)
```

The `stringsAsFactors = TRUE` parameter tells R to read any character or text variables as classes or categories and not as just text.

Recall what we should do to familiar ourselves with the dataset? We need to now how many rows and columns of the dataset, what are the variables, what types of these variables, and we may want to view the dataset for a quick scan?

Therefore, we need these functions: `dim()` or `ncol()` and `nrow()`, `names()`, `str()`, and `View()`.


```
# know how many rows and columns
dim(df)

# know the names of the variables
names(df)

# know types and examples of these variables
str(df)

# open a view window to scan the dataset
View(df)
```

Question 1. How many people have been included in this survey?

In Week 2, we also know a simple way to help us get a quick look at the descriptive summary of all the variables:

```
# summary all variables
summary(df)
```

Okay, if for now you get yourself familiar with the dataset you are going to work on, let's move on.

3.4 Sampling of numerical variables

As discussed in the lecture, the **Standard Error (SE)** of the sample mean indicates how much the estimated mean is expected to vary from sample to sample. Together, the sample mean and SE help us assess how accurately our sample reflects the true population mean. A smaller SE indicates our sample estimate is likely closer to the actual population value.

From the steps of familiarizing with the dataset, you should already know that there are 33,847 people included in this dataset. Recall our lecture, these 33,847 people should be selected by some sampling methods - random or stratified or multistaged - but in any case, they are a sample of the entire world population.

Here, in this practical, let's assume that our survey data is about a certain population. With such assumption, our population is the people what completed the survey. Thus in this section, we are pretending that the survey respondent are the population.

We can test the mean and standard deviation of the `hh_income_net` of the population using the functions `mean()` and `sd()`. In Week 2, we have learnt that the mean tells you the center of your numerical variable, and the standard deviation reveals how spread out the variable are around the mean.

```
# calculate the mean value of the variable  
mean(df$hh_income_net)
```

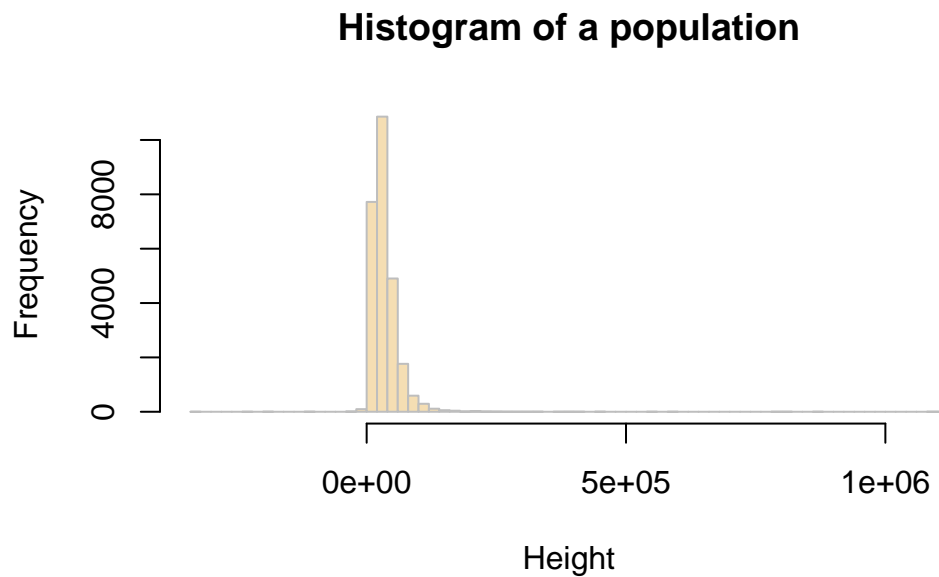
```
[1] 34485.51
```

```
# calculate the standard deviation of the variable  
sd(df$hh_income_net)
```

```
[1] 29379.1
```

You can examine the `income_net` distribution by creating a histogram as using the code below (you should be very competent of this from Week 2):

```
# plot a histogram for the variable income_net in df  
hist(df$hh_income_net, breaks = 100,  
     main = "Histogram of a population",  
     xlab = "Height",  
     col = "wheat",  
     border = "grey")
```



3.4.1 Sample mean and Standard Error

In the example below, we will work with the `hh_income_net` variable and an initial sample of 10 people (observations).

You can assess a sample of 10 observations (survey respondents) taken from the population at random using the `sample` function.

```
# Create a sample of the population
sample_10 <- slice_sample(df, n=10)
```

```
# sample mean value
mean(sample_10$hh_income_net)
```

```
[1] 38677.6
```

```
# sample standard deviation
sd(sample_10$hh_income_net)
```

```
[1] 29043.46
```

You will get a different values for these to the ones your friends get (and the ones created below), because you will have each extracted a different **sample** from the population - this is **random** sampling!

You should calculate the **Standard Error (SE)** of the sample mean by:

$$SE(\bar{x}) = \frac{SD}{\sqrt{n}}$$

In R we run the code below:

```
# sample SE
sd(sample_10$hh_income_net)/sqrt(length(sample_10$hh_income_net))
```

```
[1] 9184.347
```

```
# or
sd(sample_10$hh_income_net)/sqrt(nrow(sample_10))
```

```
[1] 9184.347
```

Question 2. Compare population mean, sample mean and SE mean from your results?

Now, to see how using a larger sample improves the estimate, you will now repeat the previous example, generating a sample, determining the sample means and SEs, but using different sample sizes: 10, 50, 100, 200.

Complete the below table for the comparison:

| Random sample size | Sample mean | Standard Error (SE) mean |
|--------------------|-------------|--------------------------|
| 10 | | |
| 50 | | |
| 100 | | |
| 200 | | |

The above should have given you the general idea that as sample size increases, the sample estimate of a population mean becomes more reliable. To show this clearly, the code below generates a trend line plot showing the impact of sample size on Standard Errors. As the sample gets larger the SE gets smaller as the sample more closely represents the true population.

```
# create a vector of sample sizes
X = seq(10, 200, 10)
# check the result
X
```

```
[1] 10 20 30 40 50 60 70 80 90 100 110 120 130 140 150 160 170 180 190
[20] 200
```

```
# create a vector of sample errors from these
# an empty vector to be populated in the for loop below
SE = vector()
# now loop through each value in X
for (i in X){
  # create a sample for the ith value in X
  # set the seed - see the Info box below
  set.seed(12)
```

```

sample.i = slice_sample(df, n=i)
# calculate the SE
se.i = sd(sample.i$hh_income_net)/sqrt(length(sample.i$hh_income_net))
# add to the SE vector
SE = append(SE, se.i)
}
# check the result
SE

```

```

[1] 11545.970  6139.880  4638.476  3893.327  3353.118  3067.783  2687.440
[8]  2391.198  2261.770  2103.450  2114.780  1979.338  1915.625  1882.336
[15]  1979.165  1878.847  1810.918  1725.824  1776.093  1705.949

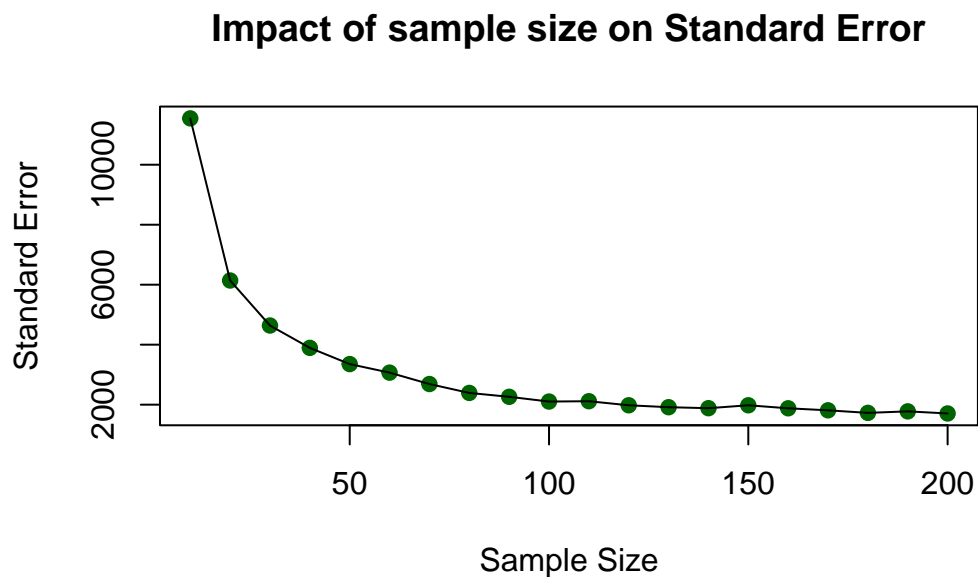
```

You can then these plot these:

```

# plot the SEs
plot(x = X, y = SE,
     pch = 19, col = "darkgreen",
     xlab = "Sample Size", ylab = "Standard Error",
     main = "Impact of sample size on Standard Error")
lines(x = X, y = SE)

```



So...as sample size increases, the sample estimate of a population mean becomes more reliable. However, you should note that this was based on one-off sample generation and looking at the SE. In the next section you will apply confidence intervals to the sample in order to assess the robustness of the sample.

3.4.2 Confidence Intervals

The idea of a **confidence interval**, CI, is a natural extension of the standard error. It allows us to define a level of confidence in our population parameter estimate gleaned from a sample:

$$\text{Sample Mean} \pm 1.96 \times SE \text{ for a 95\% interval}$$

We can use the `qnorm` function. `qnorm(p)` gives the value on the normal distribution such that the cumulative probability up to that value equals `p`. It is used at here to calculate the errors around the sample mean under an assumption of a normal distribution of the population (hence the *norm* bit of `qnorm`):

```
# you have already created the sample with
# sample_10 <- slice_sample(df, n=10)
m <- mean(sample_10$hh_income_net)
m
std <- sd(sample_10$hh_income_net)
n <- length(sample_10$hh_income_net)

error <- qnorm(0.975)*std/sqrt(n)

lower.bound <- m-error
upper.bound <- m+error
```

```
# upper and lower bound
upper.bound
lower.bound
```

This is the 95% confidence interval, for the rest 2.5% in the lower tail and 2.5% in the upper tail. This is why you may find in the code, we use `qnorm(0.975)` as that refers to the cutoff points in two directions). Again, you may find your values subtly different due to random sample. But in my case, **the mean household net income is 42,110 pound, and we are 95% confident that the mean household net income is between 22,842 pound and 61,376 pound**. This can be a really fuzzy estimates, and it is because we only use 10 respondents from the survey and by using such as small sample size, the CI window would be very fuzzy.

Let's change the above code to a sample size of 1000:

```
# change sample size from 10 to 1000 this time
sample_1k <- slice_sample(df, n=1000)
m <- mean(sample_1k$hh_income_net)
m
std <- sd(sample_1k$hh_income_net)
n <- length(sample_1k$hh_income_net)

error <- qnorm(0.975)*std/sqrt(n)

lower.bound <- m-error
upper.bound <- m+error

lower.bound
upper.bound
```

Now you see that the sample mean has changed to 36,680, which indicates that **the mean household net income is 36,680 pound, and we are 95% confident that the mean household net income is between 33,934 pound and 39,427 pound**. Therefore, with this sample size as 1000 sample, the CI window is much narrow and the result is more solid.

You may try out if we use all the respondents as the sample, what is the result? In this case, we don't need to `slice_sample()`, and the `n` should be `n = nrow(df)`.

3.5 Sampling of categorical variable

Now, we switch our focus to the categorical variable `happiness` in the dataframe. Similar to what we have done in last week, we first wish to have a descriptive summary of the categories and frequency of categories for categorical variables. We can use `table()` here:

```
# categories and frequencies
table(df$happiness)
```

| | | | |
|--------------|----------------|------------|--------------|
| Fairly happy | Fairly unhappy | Very happy | Very unhappy |
| 8346 | 1552 | 15548 | 1079 |

You may also remember how to create a table to display the frequency and proportion of each category:

```
# use df to create a table, including frequency and percentage of each category
df %>%
  count(happiness) %>%
  mutate(pct = round(n / sum(n) * 100,1))
```

| | happiness | n | pct |
|---|----------------|-------|------|
| 1 | Fairly happy | 8346 | 31.5 |
| 2 | Fairly unhappy | 1552 | 5.9 |
| 3 | Very happy | 15548 | 58.6 |
| 4 | Very unhappy | 1079 | 4.1 |

The `pct` in the table, is the percentage of respondents in our Family Resource Survey sample who reported each level of happiness. In other words, they are the **sample proportion** for variable happiness in this dataset.

Meanwhile, recall from Week 2 lecture, there are two types of categorical variables: **nominal** and **ordinal**. The key difference between them is whether the categories have a nature ordering. In this case, the happiness variable is **ordinal**, because its categories follow a meaningful order: *Very unhappy*, *Fairly unhappy*, *Fairly happy*, and *Very happy*.

Therefore, as we wish to display the variable with an order, we use R function `factor()` to do so:

```
#set the ordering of categories
df$happiness <- factor(
  df$happiness,
  levels = c("Very unhappy",
             "Fairly unhappy",
             "Fairly happy",
             "Very happy"),
  ordered = TRUE
)
```

If we now re-run the code to generate the frequency table and proportions, you will notice that the categories appear in a more logical order, reflecting the ordinal structure of the variable.

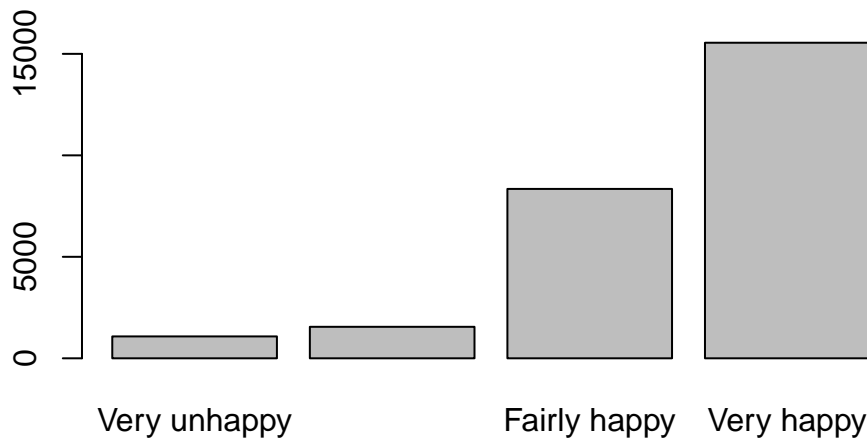
This time, we label the proportion column as \hat{p} (**p-hat**) to emphasise that these values represent the **sample proportions** for each happiness category.

```
df %>%
  count(happiness) %>%
  mutate(p_hat = round(n / sum(n) * 100,1))
```


| | happiness | n | p_hat |
|---|----------------|-------|-------|
| 1 | Very unhappy | 1079 | 4.1 |
| 2 | Fairly unhappy | 1552 | 5.9 |
| 3 | Fairly happy | 8346 | 31.5 |
| 4 | Very happy | 15548 | 58.6 |

We use `barplot()` from last week to present the happiness variable distribution in the Family Resource Survey:

```
#use bar plot to visualise the distribution of categorical variable
tab=table(df$happiness)
barplot(tab)
```



Question 3. What pattern you can identify from the barplot about happiness in this sample?

3.5.1 Sample proportion and Standard Error of the proportion

In this section, we will focus on the **Very happy** category, we can check the sample mean of Very happy from the eariler work and the `p_hat` is 58.6%. This can also be calculated by `mean()`:

```
#calculate sample mean
p_hat = mean(df$happiness == "Very happy")
p_hat
```

```
[1] 0.586164
```

Now let's calculate the **Standard Error**. Recap the equation of Standard Error for a proportion:

$$SE(\hat{p}) = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

So, we need the sample size **n** first and then use **n** and **p_hat** to calculate the **SE_p**. If we select all the respondents in the dataset as the sample, then the n is equal to the **nrow()**.

```
#calculate Standard Error
n = nrow(df)
n
```

```
[1] 26525
```

```
SE_p <- sqrt(p_hat * (1 - p_hat) / n)
SE_p
```

```
[1] 0.003024099
```

Therefore, now we get the interpretation that the “Very happy” sample proportion typically varies by about 0.003 or 0.30 percentage.

We can also test **n** from 10, 100, 1000 to all to plot the **SE_p** and see how sample size increased will reduce the Standard Error of proportion:

```
X = c(10, 100, 1000, 2000, 5000, 10000, n)
SE_p = vector()
# now loop though each value in X
for (i in X){
  # create a sample for the ith value in X
  # set the seed - see the Info box below
  set.seed(20)
```

```

sample.i = slice_sample(df, n=i)
# calculate the sample mean of proportion
p_hat = mean(sample.i$happiness == "Very happy")
# calculate the SE of proportion
se.i = sqrt(p_hat * (1 - p_hat) / i)
# add to the SE vector
SE_p = append(SE_p, se.i)
}
# check the result
SE_p

```

```

[1] 0.158113883 0.049355851 0.015409607 0.010928626 0.006966181 0.004908920
[7] 0.003024099

```

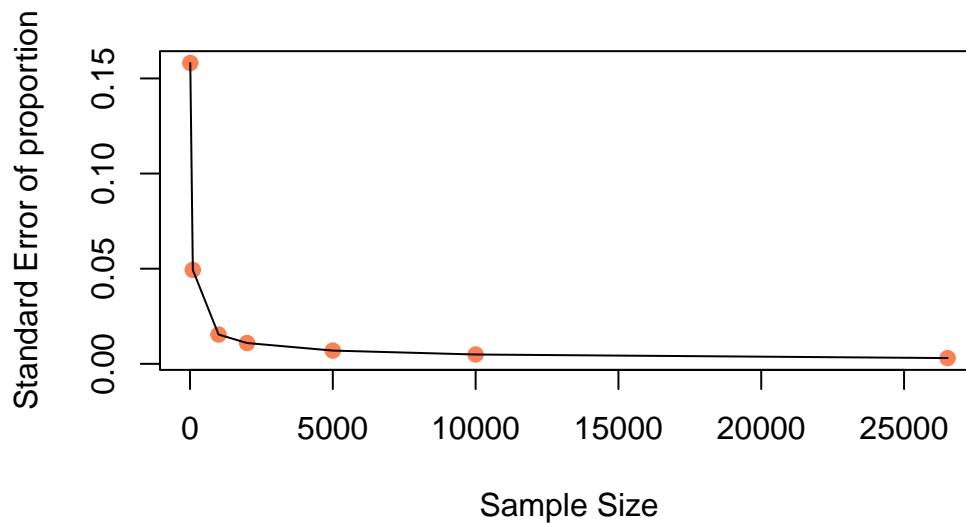
We can then plot these:

```

# plot the SEs
plot(x = X, y = SE_p,
     pch = 19, col = "coral",
     xlab = "Sample Size", ylab = "Standard Error of proportion",
     main = "Impact of sample size on Standard Error of the proportion Very happy")
lines(x = X, y = SE_p)

```

Impact of sample size on Standard Error of the proportion Very



Again, we know that with the increase of sample size n , the Standard Error decreases.

3.5.2 Confidence Intervals for a proportion (95%)

We now calculate a **95% confidence interval** for the proportion of respondents who reported being *Very happy* in our sample.

To do this, we use the full sample size n . Recall that for a proportion, CI for a proportion (95%) is

- Sample proportion \pm (critical value \times standard error)

The critical value is again the `qnorm(0.975)` as we did for the numerical variable, if you `qnorm(0.975)` then you may find it is 1.96. The below R code help use to calculate the CI windows of sample mean, when sample size is n .

```
#CI windows
p_hat = mean(df$happiness == "Very happy")
p_hat
```

```
[1] 0.586164
```

```
n <- nrow(df)
SE_p <- sqrt(p_hat * (1 - p_hat) / n)
error_p <- qnorm(0.975) * SE_p
lower.bound <- p_hat - error_p
upper.bound <- p_hat + error_p
lower.bound
```

```
[1] 0.5802369
```

```
upper.bound
```

```
[1] 0.5920911
```

Therefore, the confidence interval can be interpreted as follows: in the current UK *Family Resources Survey*, the sample proportion of respondents reporting being *Very happy* is 58.6%. The 95% confidence interval indicates that we are 95% confident that the true population proportion lies between 58.0% and 59.2%.

3.6 Compare to World Value Survey

We know that when the sample changes, the survey results may also differ. Now we turn to a different dataset — the *World Values Survey* — and repeat the same analytical process to examine how the happiness results compare.

3.6.1 Happiness in different survey sample

As usual, we first load the World Value Survey data:

```
# use read.csv to load a CSV file
dat <- read.csv("world_value_survey.csv", stringsAsFactors = TRUE)
```

Again, we wish to treat the happiness variable with its natural ordering:

```
#set ordering
dat$happiness <- factor(
  dat$happiness,
  levels = c("Very unhappy",
             "Fairly unhappy",
             "Fairly happy",
             "Very happy"),
  ordered = TRUE
)
```

We also want to have a descriptive summary of the count in each category and compute the percentage distribution:

```
# descriptive summary table
dat %>%
  count(happiness) %>%
  mutate(p_hat2 = round(n / sum(n) * 100, 1))
```

| | happiness | n | p_hat2 |
|---|----------------|-------|--------|
| 1 | Very unhappy | 2517 | 1.6 |
| 2 | Fairly unhappy | 18032 | 11.8 |
| 3 | Fairly happy | 87342 | 57.1 |
| 4 | Very happy | 45118 | 29.5 |

Not just as the table above, a barplot would better support the descriptive summary visually:

```
# barplot visualise
tab = table(dat$happiness)
tab
```

| Very unhappy | Fairly unhappy | Fairly happy | Very happy |
|--------------|----------------|--------------|------------|
| 2517 | 18032 | 87342 | 45118 |

```
barplot(tab)
```



As what we have done earlier, we want to focus on the Very happy category so we first calculate the sample proportion of respondents who are Very happy.

```
#sample proportion of Very happy
p_hat2 = mean(dat$happiness == "Very happy")
p_hat2
```

```
[1] 0.2948715
```

To estimate the Standard Error, we use the same coding, but this time, the sample size is called `k` and all the respondents are included:

```
#Standard Error for the proportion of Very happy
k = nrow(dat)
k
```

```
[1] 153009
```

```
SE_p2 <- sqrt(p_hat2 * (1 - p_hat2) / k)
SE_p2
```

```
[1] 0.001165714
```

The estimated proportion of respondents who are *Very happy* typically varies by about 0.0012 (or 0.12 percentage points) from sample to sample due to random sampling variation.

And the Confidence Intervals:

```
p_hat2
```

```
[1] 0.2948715
```

```
SE_p2
```

```
[1] 0.001165714
```

```
error_p2 <- qnorm(0.975) * SE_p2
lower.bound2 <- p_hat2 - error_p2
upper.bound2 <- p_hat2 + error_p2
lower.bound2
```

```
[1] 0.2925868
```

```
upper.bound2
```

```
[1] 0.2971563
```

The sample proportion of respondents who are *Very happy* is estimated to be around 29.5%. The 95% confidence interval suggests that we are 95% confident that the true population proportion lies between 29.3% and 29.7%.

Compare to our earlier results by using UK Family Resource Survey, the findings are quite lower. You may have your inference that this is because earlier respondents are all from the UK but now it is from all over the world. Therefore, let's filter only UK respondents from this World Value Survey:

```
# filter UK respondents from the whole survey sample
dat_uk = dat %>% filter(english_short_name == "United Kingdom")
```

```
#calculate sample mean proportion for UK sample only
dat_uk %>%
  count(happiness) %>%
  mutate(p_hat3 = round(n / sum(n) * 100,1))
```

| | happiness | n | p_hat3 |
|---|----------------|------|--------|
| 1 | Very unhappy | 34 | 0.7 |
| 2 | Fairly unhappy | 335 | 7.1 |
| 3 | Fairly happy | 2572 | 54.6 |
| 4 | Very happy | 1772 | 37.6 |

```
p_hat3 = mean(dat_uk$happiness == "Very happy")
p_hat3
```

```
[1] 0.3759813
```

Question 4. Compute the SE and CI for the UK sample only dataframe `dat_uk`. How you will interpret the results?

3.6.2 Between regions and countries

You may have noticed that, even within the World Values Survey, the sample proportion of respondents reporting being *Very happy* in the UK is higher than the overall sample proportion.

We can now continue to explore how happiness varies across different regions and countries to see whether similar patterns emerge elsewhere.

To fulfill this task, we shall use the cross-tabulation method we did last week, by using `table()` between `region_name` and `happiness`:


```
#cross-tabulate region vs happiness
table(dat$region_name, dat$happiness)
```

| | Very unhappy | Fairly unhappy | Fairly happy | Very happy |
|----------|--------------|----------------|--------------|------------|
| Africa | 421 | 1627 | 4829 | 2645 |
| Americas | 287 | 2868 | 11541 | 9053 |
| Asia | 866 | 5468 | 26994 | 14992 |
| Europe | 927 | 7894 | 42287 | 17539 |
| Oceania | 16 | 175 | 1691 | 889 |

You may think, okay, but it is very hard to draw conclusions based on the counts. You are right. Let's change the count table into a proportion one by using function `prop.table()`.

To know the function, we can first ask R's help:

```
?prop.table
```

On your RStudio's right-hand Help pane, we can learn that we can set parameter `margin` in to 1 to indicate the division should be done by row (each row sums up to as 1), and `margin = 2` as the division by columns (each column sums up to 1).

Let's try both `margin=1` and `margin=2`:

```
#get the cross-tabulation result and present in proportion format, margin = 1
tab_region_happy = table(dat$region_name, dat$happiness)
prop.table(tab_region_happy, margin = 1) * 100
```

| | Very unhappy | Fairly unhappy | Fairly happy | Very happy |
|----------|--------------|----------------|--------------|------------|
| Africa | 4.4213401 | 17.0867465 | 50.7141357 | 27.7777778 |
| Americas | 1.2084719 | 12.0762979 | 48.5957303 | 38.1194998 |
| Asia | 1.7922185 | 11.3162252 | 55.8650662 | 31.0264901 |
| Europe | 1.3503868 | 11.4994100 | 61.6006526 | 25.5495506 |
| Oceania | 0.5774089 | 6.3154096 | 61.0249008 | 32.0822808 |

```
round(prop.table(tab_region_happy, margin = 1) * 100, 2) #round(...,2) to keep only 2 digits
```

| | Very unhappy | Fairly unhappy | Fairly happy | Very happy |
|----------|--------------|----------------|--------------|------------|
| Africa | 4.42 | 17.09 | 50.71 | 27.78 |
| Americas | 1.21 | 12.08 | 48.60 | 38.12 |
| Asia | 1.79 | 11.32 | 55.87 | 31.03 |
| Europe | 1.35 | 11.50 | 61.60 | 25.55 |
| Oceania | 0.58 | 6.32 | 61.02 | 32.08 |

A vague calculation you may notice that each row (region) sums to 100 (as in the above codes, we multiplied the proportion by 100 to express them as percentage). This means the table presents **row percentages**: within each region, the percentages show how respondents are distributed across the different levels of happiness. In other words, the table allows us to compare how the distribution of happiness varies across regions within the sample.

If we simply change `margin = 1` to `margin = 2`, the results will be different because we are now calculating column percentages instead of row percentages.

```
round(prop.table(tab_region_happy, margin = 2) * 100, 2)
```

| | Very unhappy | Fairly unhappy | Fairly happy | Very happy |
|----------|--------------|----------------|--------------|------------|
| Africa | 16.73 | 9.02 | 5.53 | 5.86 |
| Americas | 11.40 | 15.91 | 13.21 | 20.07 |
| Asia | 34.41 | 30.32 | 30.91 | 33.23 |
| Europe | 36.83 | 43.78 | 48.42 | 38.87 |
| Oceania | 0.64 | 0.97 | 1.94 | 1.97 |

With `margin = 2`, each column sums to 100. This means we are looking at the distribution of regions within each happiness category. In other words, instead of asking: “Within each region, how is happiness distributed?” we are now asking: “Among respondents at each happiness level, how are they distributed across regions?” - This changes the interpretation entirely.

Now it's the time to check out the comparison among countries. Here we use variable `english_short_name` as the country to cross-tabulate `happiness`:

```
tab_country_happy = table(dat$english_short_name, dat$happiness)
round(prop.table(tab_country_happy, margin = 1) * 100, 1) #round(...,1) to keep only 2 digits
```

Question 5. What conclusion you can draw from the findings? Change margin from 1 to 2, what new information you can learn from the results?

Again, in the following section, we want focus on the proportion of very happy in each country only, we can use two functions `group_by()` `%>% summarise()` to fulfill this request:

```
#group dat by the country name, in each group, calculate the mean proportion of happiness rep
df_ctype_very_happy <- dat %>%
  group_by(english_short_name) %>%
  summarise(p_hat = mean(happiness == "Very happy", na.rm = TRUE)) #if any rows include NA v
```

```
df_ctype_very_happy
```

This produces one row per country, with the sample proportion of respondents who are *Very happy*.

We can get the top 1 country by using function `which.max()`. This function returns the row number corresponding to the maximum value of `p_hat`. We can then use that row number to extract the full row from the dataframe: to get the row number of that country, and then ask R to return that row for us:

```
which.max(df_ctype_very_happy$p_hat)
```

```
[1] 44
```

```
df_ctype_very_happy[which.max(df_ctype_very_happy$p_hat),]
```

```
# A tibble: 1 x 2
  english_short_name p_hat
  <fct>             <dbl>
1 Kyrgyz Republic  0.677
```

To identify the top 10 happiest countries (based on the proportion of respondents who are *Very happy*), we can use `arrange()` function to sort the `df_ctype_very_happy` dataframe by `p_hat`.

If we place a minus sign - before `p_hat` (i.e. `arrange(-p_hat)`), the dataframe will be sorted in **decreasing order**, from highest to lowest proportion. Without the minus sign, the sorting will be in **ascending order**, from lowest to highest.

You can also type `?arrange` in RStudio to view the official documentation and see additional details about how the function works.

Now, top 10 Very happy:

```
df_ctype_very_happy %>% arrange(-p_hat) %>% head(10)
```

```
# A tibble: 10 x 2
  english_short_name p_hat
  <fct>             <dbl>
1 Kyrgyz Republic  0.677
2 Tajikistan        0.626
3 Ecuador           0.618
4 Mexico            0.588
5 Colombia          0.574
6 Uzbekistan        0.519
7 Guatemala         0.511
8 Philippines        0.511
9 Puerto Rico       0.510
10 Kenya           0.509
```

bottom 10 Very happy:

```
df_etry_very_happy %>% arrange(p_hat) %>% head(10)
```

```
# A tibble: 10 x 2
  english_short_name p_hat
  <fct>             <dbl>
1 South Korea       0.0410
2 Egypt            0.0644
3 Lithuania         0.0918
4 Hong Kong         0.111
5 Morocco          0.126
6 Latvia           0.132
7 Slovakia          0.133
8 Macao            0.145
9 Russia           0.149
10 Greece           0.154
```

3.6.3 Map happiness around the world

As usual, let's finish today's practical by a map made by ourselves.

First, load the world map boundaries from your working folder. Before producing the final map, we can generate a quick check map using `qtm()` to ensure the spatial data has loaded correctly. We also set `tmap` to interactive mode using `tmap_mode("view")`, so the map can be explored dynamically (e.g., zooming and clicking on countries) before creating the final visualisation.

```
#load the boundary dataset
world_map <- st_read("world_map.geojson")
```

```
Reading layer `world_map' from data source
  `C:\Users\ziye\Documents\GitHub\quant\labs\world_map.geojson'
  using driver `GeoJSON'
Simple feature collection with 203 features and 14 fields
Geometry type: MULTIPOLYGON
Dimension:      XY
Bounding box:   xmin: -180 ymin: -58.49861 xmax: 180 ymax: 83.6236
Geodetic CRS:   WGS 84
```

```
#set as interactive map
tmap_mode("view")
#quick mapping
qtm(world_map)
```

Click on any country in the world map within the Viewer pane on the right-hand side of RStudio. You will notice that the variable `english_short` contains the country name.

We can use this variable as the key to join the spatial map data with our country-level happiness dataframe, since it corresponds to the `english_short_name` variable in our survey dataset.

Before join them, let's first inspect the structure of both dataframe:

```
dim(world_map)
```

```
[1] 203  15
```

```
dim(df_ctry_very_happy)
```

```
[1] 91  2
```

We can learn that the world map has 203 countries, but our `df_ctry_very_happy` only have 91, this means that not all the countries from the `world_map` will be able to plot based on the `p_hat` from `df_ctry_very_happy`. As a result, some countries will have missing values (NA) for the happiness proportion and therefore will not be coloured according to `p_hat` in the final map.

However, this is perfectly acceptable — those countries will simply appear with the default missing-value styling, while the countries included in our survey data will be shaded according to their estimated proportion of respondents who are Very happy.


Let's use function `left_join()` to join the two dataframes:

```
#join our newly created dataframe df_ctry_very_happy to the world map by the shared columns
world_happiness = left_join(world_map, df_ctry_very_happy, by=c("english_short"="english_shor
```

A short code sentence would serve your first try on the new dataframe `world_happiness`.

Don't forget to click the  from your RStudio right-bottom Viewer pane for a full screen interactive map in your explorer:

```
#View the map as an interview map
tm_shape(world_happiness) +
  tm_polygons("p_hat")
```

To make it looking better, let's change the color palette and breaks. Click the  from your RStudio right-bottom Viewer pane for a full screen interactive map - browse the map and interact with it to find how UK compare to other countries on this map? Any other findings interest you?

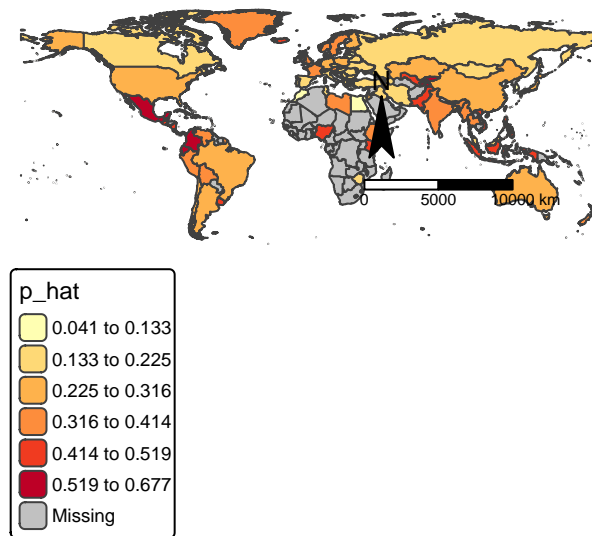
```
#plot the map as a static map, with color palette as Yellow-Orange_Red, with break styles as

map = tm_shape(world_happiness) +
  tm_polygons("p_hat",
              fill.scale = tm_scale(values="YlOrRd",
                                   style = "jenks",n=6))+
  tm_layout(main.title = "Mean proportion of Very happy in the World Value Survey",
            main.title.size=1.2,
            frame = FALSE) +
  tm_compass(position = c("right", "bottom")) +
  tm_scalebar(position = c("right", "bottom"))
map
```

At last, we can use `tmap_mode("plot")` to change the interactive map into a static one. Also we can save the map on your disk:

```
tmap_mode("plot")
map
tmap_save(map, "Week3 map.png")
```

Mean proportion of Very happy in the World Value Survey



Check in your folder and you shall find the map made by yourself.

3.7 Formative Tasks

Task 1 Use the `income_gross` variable from the UK Family Resource Survey, compare how increase sample size from 5000 to 20000 would affect the sample mean, standard error and CI of the `income_gross`. How to interpret the results?

```
n1 = 5000
n2 = 20000

# sample by size
sample_5k = slice_sample(df, n=n1)
sample_20k = slice_sample(df, n=n2)

#mean
m.5k = mean(sample_5k$income_gross)
m.20k = mean(sample_20k$income_gross)

m.5k
```

```
[1] 22224.9
```

```
m.20k
```

```
[1] 22076.63
```

```
#standard deviation
std.5k = sd(sample_5k$income_gross)
std.20k = sd(sample_20k$income_gross)

#standard error
se.5k = std.5k/sqrt(n1)
se.20k = std.20k/sqrt(n2)

se.5k
```

```
[1] 434.672
```

```
se.20k
```

```
[1] 174.1159
```

```
#CI
error.5k <- qnorm(0.975)*se.5k
error.20k <- qnorm(0.975)*se.20k

lower.bound.5k <- m.5k-error.5k
upper.bound.5k <- m.5k+error.5k

lower.bound.20k <- m.20k-error.20k
upper.bound.20k <- m.20k+error.20k

lower.bound.5k
```

```
[1] 21372.96
```

```
upper.bound.5k
```

```
[1] 23076.85
```



```
lower.bound.20k
```

```
[1] 21735.37
```

```
upper.bound.20k
```

```
[1] 22417.89
```

Task 2 Use the `health` variable from the UK Family Resource Survey, what is the sample mean proportion of **Very Good** health in the whole sample? The standard error and CI? How to interpret the results?

```
n = nrow(df)
p_hat = mean(df$health == "Very Good")
p_hat
```

```
[1] 0.2820358
```

```
SE_p <- sqrt(p_hat * (1 - p_hat) / n)
error_p <- qnorm(0.975) * SE_p
lower.bound <- p_hat - error_p
upper.bound <- p_hat + error_p
lower.bound
```

```
[1] 0.2766205
```

```
upper.bound
```

```
[1] 0.2874511
```

Task 3 Use the `health` variable from the World Value Survey, compare your the sample mean proportion of **Very good** health with your findings in Task 2.

```
n = nrow(dat)
p_hat = mean(dat$health == "Very good")
p_hat
```

```
[1] 0.2204315
```

```
SE_p <- sqrt(p_hat * (1 - p_hat) / n)
error_p <- qnorm(0.975) * SE_p
lower_bound <- p_hat - error_p
upper_bound <- p_hat + error_p
lower_bound
```

```
[1] 0.2183544
```

```
upper_bound
```

```
[1] 0.2225086
```

Task 4 Use the `trust_strangers` variable from the World Value Survey, compare between different regions (variable use `sub_region_name`). Which region more likely to trust strangers completely?

```
tab = table(dat$sub_region_name, dat$trust_strangers)
round(prop.table(tab, margin = 1) * 100, 1)
```

| | Do not trust at all | Do not trust much |
|---------------------------------|---------------------|-------------------|
| Australia and New Zealand | 8.2 | 36.5 |
| Central Asia | 41.5 | 41.5 |
| Eastern Asia | 19.0 | 60.6 |
| Eastern Europe | 24.0 | 48.6 |
| Latin America and the Caribbean | 48.1 | 35.3 |
| Northern Africa | 31.2 | 43.1 |
| Northern America | 11.6 | 42.2 |
| Northern Europe | 8.7 | 32.7 |
| South-eastern Asia | 29.3 | 49.2 |
| Southern Asia | 38.3 | 32.4 |
| Southern Europe | 27.6 | 46.1 |
| Sub-Saharan Africa | 33.9 | 34.7 |
| Western Asia | 33.8 | 45.4 |
| Western Europe | 10.8 | 38.1 |
| | Trust completely | Trust somewhat |
| Australia and New Zealand | 1.2 | 54.1 |
| Central Asia | 2.4 | 14.6 |
| Eastern Asia | 1.8 | 18.7 |

| | | |
|---------------------------------|-----|------|
| Eastern Europe | 2.5 | 24.9 |
| Latin America and the Caribbean | 1.9 | 14.7 |
| Northern Africa | 2.3 | 23.4 |
| Northern America | 2.0 | 44.2 |
| Northern Europe | 5.2 | 53.5 |
| South-eastern Asia | 3.0 | 18.5 |
| Southern Asia | 5.5 | 23.7 |
| Southern Europe | 2.1 | 24.2 |
| Sub-Saharan Africa | 4.8 | 26.6 |
| Western Asia | 2.3 | 18.4 |
| Western Europe | 3.1 | 48.0 |

```
result = round(prop.table(tab, margin = 1) * 100, 1) %>% as.data.frame.matrix()
result
```

| | Do not trust at all | Do not trust much |
|---------------------------------|---------------------|-------------------|
| Australia and New Zealand | 8.2 | 36.5 |
| Central Asia | 41.5 | 41.5 |
| Eastern Asia | 19.0 | 60.6 |
| Eastern Europe | 24.0 | 48.6 |
| Latin America and the Caribbean | 48.1 | 35.3 |
| Northern Africa | 31.2 | 43.1 |
| Northern America | 11.6 | 42.2 |
| Northern Europe | 8.7 | 32.7 |
| South-eastern Asia | 29.3 | 49.2 |
| Southern Asia | 38.3 | 32.4 |
| Southern Europe | 27.6 | 46.1 |
| Sub-Saharan Africa | 33.9 | 34.7 |
| Western Asia | 33.8 | 45.4 |
| Western Europe | 10.8 | 38.1 |
| | Trust completely | Trust somewhat |
| Australia and New Zealand | 1.2 | 54.1 |
| Central Asia | 2.4 | 14.6 |
| Eastern Asia | 1.8 | 18.7 |
| Eastern Europe | 2.5 | 24.9 |
| Latin America and the Caribbean | 1.9 | 14.7 |
| Northern Africa | 2.3 | 23.4 |
| Northern America | 2.0 | 44.2 |
| Northern Europe | 5.2 | 53.5 |
| South-eastern Asia | 3.0 | 18.5 |
| Southern Asia | 5.5 | 23.7 |
| Southern Europe | 2.1 | 24.2 |

| | | |
|--------------------|-----|------|
| Sub-Saharan Africa | 4.8 | 26.6 |
| Western Asia | 2.3 | 18.4 |
| Western Europe | 3.1 | 48.0 |

```
which.max(result$`Trust completely`)
```

```
[1] 10
```

```
result[which.max(result$`Trust completely`),]
```

| | Do not trust at all | Do not trust much | Trust completely |
|---------------|---------------------|-------------------|------------------|
| Southern Asia | 38.3 | 32.4 | 5.5 |
| | Trust somewhat | | |
| Southern Asia | 23.7 | | |

Task 5 Use the `immigrant_impact` variable from the World Value Survey, create a new dataframe to include the sample mean proportion of in each country/place (variable needed: `english_short_name`) who holds **Very bad** attitude of immigrant impact. Identify the five countries/places with the lowest mean proportion in the sample.

```
df_immigrant <- dat %>%
  group_by(english_short_name) %>%
  summarise(p_hat = mean(immigrant_impact == "Very bad", na.rm = TRUE)) # NA value removed

df_immigrant
```

```
# A tibble: 91 x 2
  english_short_name p_hat
  <fct>             <dbl>
1 Albania           0.0320
2 Andorra           0.00402
3 Argentina         0.0567
4 Armenia           0.0397
5 Australia         0.145
6 Austria           0.103
7 Azerbaijan        0.0643
8 Bangladesh        0.0387
9 Belarus           0.0293
10 Bolivia          0.0724
# i 81 more rows
```

```
df_immigrant %>% arrange(p_hat) %>% head(5)
```

```
# A tibble: 5 x 2
  english_short_name p_hat
  <fct>             <dbl>
1 Andorra           0.00402
2 South Korea       0.00884
3 Macao             0.00888
4 China             0.0108
5 Peru              0.0110
```

4 Lab: Correlation, data reliability and the issue of scale - Health

4.1 Overview

“No society can legitimately call itself civilised if a sick person is denied medical aid because of lack of means.”

(Aneurin Bevan, 1952, *In Place of Fear*, London: Heinemann.)

Health and wealth are closely connected.

- “The greatest wealth is health.” (Proverb)
- “The richer you are, the better your health.” (Michael Marmot, 2015, *The Guardian*, 11/09/2015.)
- Poor national health reduces workforce participation, productivity and growth, while increasing demand for health, care and welfare services (Geoffrey Filkin et al., 2024, *Health is Wealth*, London: The King’s Fund).


This practical explores the **geography of health** and its relationship to socio-economic conditions. It builds on skills developed in previous weeks and introduces additional analytical techniques.

In particular, you will:

- Link data from two different datasets
- Visualise the relationships between different types of variables
- Calculate appropriate measures of correlation for different types of variables
- Consider Modifiable Areal Unit Problem (MAUP)

4.2 Prepare your working environment

- For this Week 4 session, create a sub-folder called Week4 in your ENVS162 folder on your M-Drive. This is exactly the step we did in Weeks 1, 2, and 3.
- Download this week’s practical datasets from Canvas Week 4:
 - LAD December 2021 EW BUC.geojson

- Census_2021_Districts.csv
 - Census_2021_Regions.csv
 - Census_2021_Counties.csv
 - Census_2021_Wards.csv
 - England_Wales_SAR.csv
- Open RStudio
 - Open a new R Script for your Week 4 work, rename it as Week4.R and save it in your newly created Week 4 folder, under M drive -> ENVS162 folder.
 - Check whether there are any previous left dataframes in your RStudio in the upper-right side Environment pane. You can always use the to clear all the dataframes in your environment and make it all clean. For the same aim, you can click the icon , or you can run the below code:

```
rm(list = ls())
```

This command and also the brush icon can both clear RStudio's memory, removing any data objects that you have previously created.

4.3 Load libraries

Exactly as what we have done for the previous weeks, before starting any analysis in R, we first need to load the essential libraries as all the functions/commands are packed in these libraries. For this week, we will still rely on `tidyverse`, `tmap` and `sf`.

```
library(tidyverse)
library(sf)
library(tmap)
```

4.4 The Geography of Health (area-level)

4.4.1 Viewing the district level boundaries

The *districts* dataset is an area-level dataset, in which each row represents one district. Let's view what this data looks like.

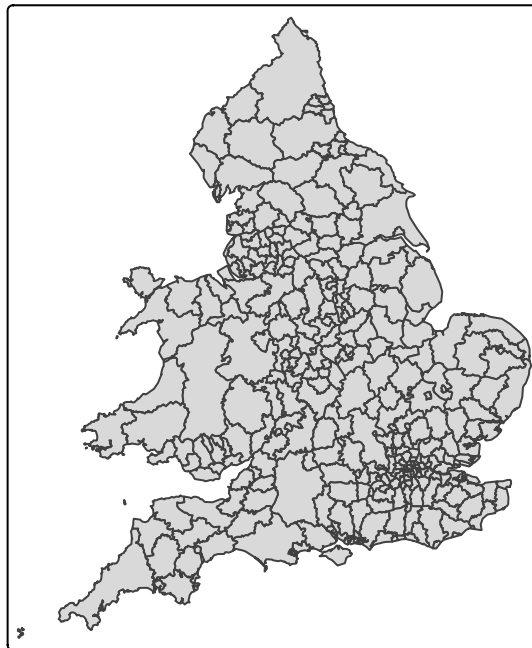
4.4.1.1 Load the district boundaries data

```
#load the boundary dataset
district_map <- st_read("LAD December 2021 EW BUC.geojson")
```

```
Reading layer `LAD December 2021 EW BUC.geojson' from data source
`C:\Users\ziye\Documents\GitHub\quant\labs\LAD December 2021 EW BUC.geojson'
using driver `GeoJSON'
Simple feature collection with 331 features and 1 field
Geometry type: MULTIPOLYGON
Dimension:      XY
Bounding box:   xmin: -6.360297 ymin: 49.88234 xmax: 1.763571 ymax: 55.8112
Geodetic CRS:   WGS 84
```

4.4.1.2 Display on a map

```
qtm(district_map)
```



These are just the boundaries of the districts across England.

4.4.1.3 View socioeconomic data at the district level

This data is contained in the `Census_2021_Districts.csv` file. Let's load this file and view its contents.

```
districts <- read.csv("Census_2021_Districts.csv", stringsAsFactors = TRUE)
```

The census dataset contains information on each district in England and Wales derived from the 2021 Census.

The dataset contains three sets of variables

- *LAD21CD* to *Country* - a set of *geographic identifiers* relating to the location of each district [There is a variable in the *district_map*]
- *Density* and *Mean_Age* - two *measures*: population density (per km²) and mean age (in years)
- *Age_0_to_19* to *Public_Transport* - a set of variables measuring the *percentage* of persons in each district that fall within the various the categories named by the variables. For example, the variable *Age_0_to_19* records the % of persons aged 0 to 19, whilst *Public_Transport* records the % of persons who commute to work by public transport.

Recall what we should do to familiar ourselves with the dataset? We need to know how many rows and columns of the dataset, what are the variables, what types of these variables, and we may want to view the dataset for a quick scan?

Therefore, we need these functions: `dim()` or `ncol()` and `nrow()`, `names()`, `str()`, and `View()`.

```
# know how many rows and columns
dim(districts)
# know the names of the variables
names(districts)
# know types and examples of these variables
str(districts)
# open a view window to scan the dataset
View(districts)
```

QUESTION 1

How many rows does the districts data have?

How many columns does the districts data have?

The *districts* dataset has a lot of data but it would be useful to visualise this data on a map.

In order to map any of the variables contained in the *districts* dataset, we need to join them to the *district_map* dataset using the `full_join()` command. We can link the two datasets using the variable they share in common, *LAD21CD*, which is the unique identifier given by the Office for National Statistics to each district captured in the 2021 Census.

```
census_map <- left_join(district_map, districts, by = "LAD21CD")
```

Note that the dataset containing the map boundaries (*district_map*) MUST be listed *before* the dataset it is being joined to (*districts*), otherwise the `left_join()` command will throw an error.

QUESTION 2

Write a line of code that tells you the number of columns in the `census_map` dataset.

```
# REPLACE THIS LINE WITH YOUR CODE THEN RUN/EXECUTE THIS BLOCK
```

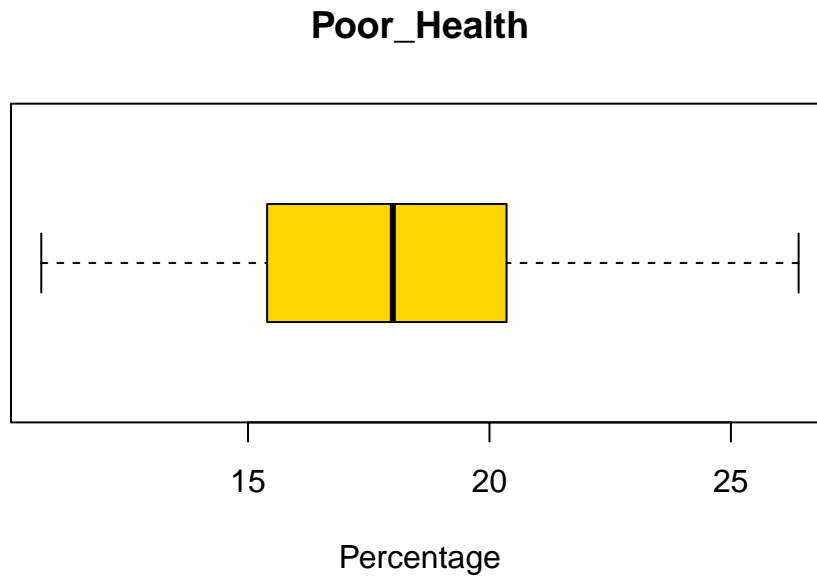
4.4.2 The spatial distribution Poor Health

One of the variables contained in the *district* dataset is *Poor_Health*, which records the % of persons in each district who self-reported their health as being either 'poor' or 'very poor'.

4.4.2.1 Examining skewness in the data

Let's now examine the distribution of this data. Recall from week 2 lab that one way to do this is using a boxplot.

```
boxplot(census_map$Poor_Health,  
        horizontal=TRUE,  
        main = "Poor_Health",  
        xlab='Percentage',  
        col = "gold")
```



QUESTION 4

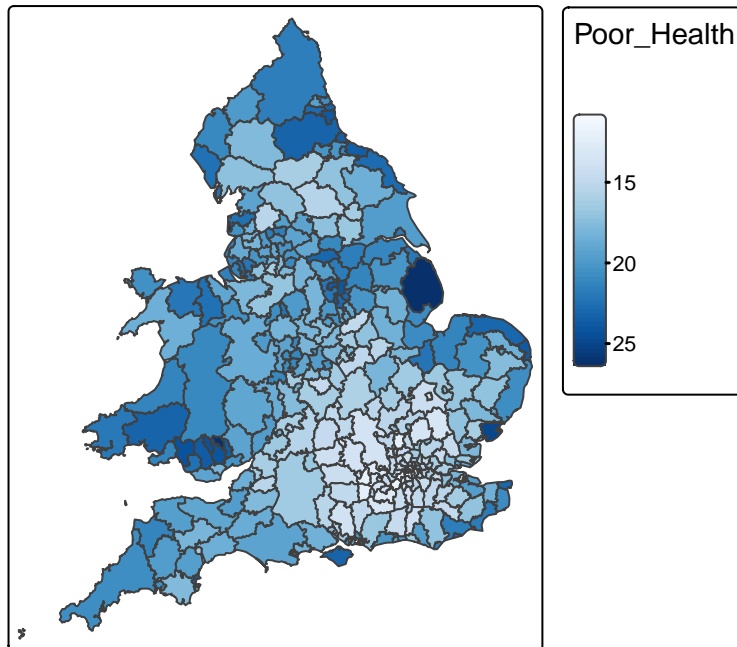
Return to the week 2 lab. Using the information provided, determine whether *Poor_Health* is skewed or not.

4.4.2.2 Map the spatial distribution of *Poor_Health*.

Use the `tm_scale_continuous()` option to make sure that areas with the highest rates of poor health have the darkest shading.

```
Poor_Health_map <-
  tm_shape(census_map) +
  tm_polygons(
    fill = "Poor_Health",
    fill.scale = tm_scale_continuous(values = "Blues")
  )

Poor_Health_map
```



Some areas have really high *Poor_Health* values. Let's now find out which district have the highest percentage of people with poor health.

```
census_map %>%
  slice(which.max(Poor_Health)) %>%
  select(`District`, Poor_Health)
```

Simple feature collection with 1 feature and 2 fields

Geometry type: MULTIPOLYGON

Dimension: XY

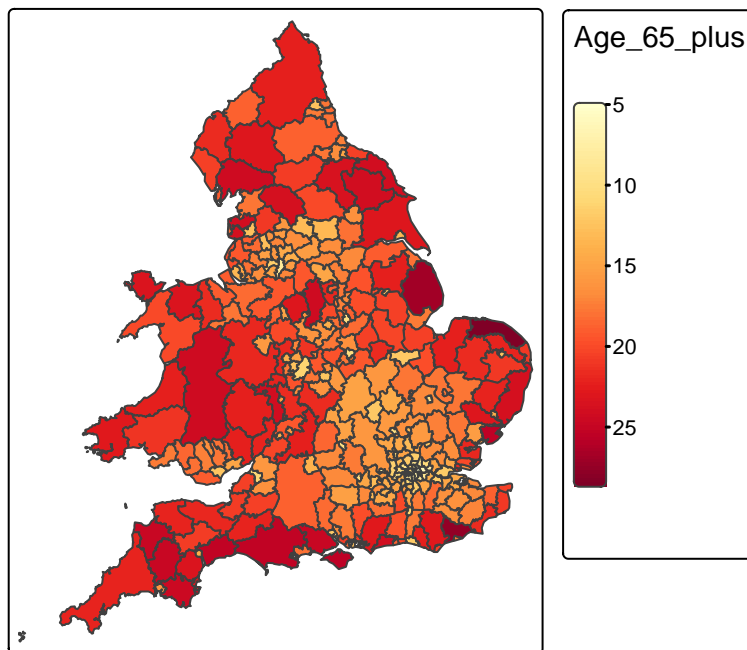
Bounding box: xmin: -0.3174608 ymin: 53.00112 xmax: 0.3556274 ymax: 53.52762

Geodetic CRS: WGS 84

| District | Poor_Health | geometry |
|----------------|-------------|--------------------------------|
| 1 East Lindsey | 26.40201 | MULTIPOLYGON (((0.2409995 5... |

4.4.2.3 Map the spatial distribution of Age_65_plus.

```
Age_65_plus_map <-  
  tm_shape(census_map) +  
  tm_polygons(  
    fill = "Age_65_plus",  
    fill.scale = tm_scale_continuous(values = "YlOrRd")  
  )  
Age_65_plus_map
```



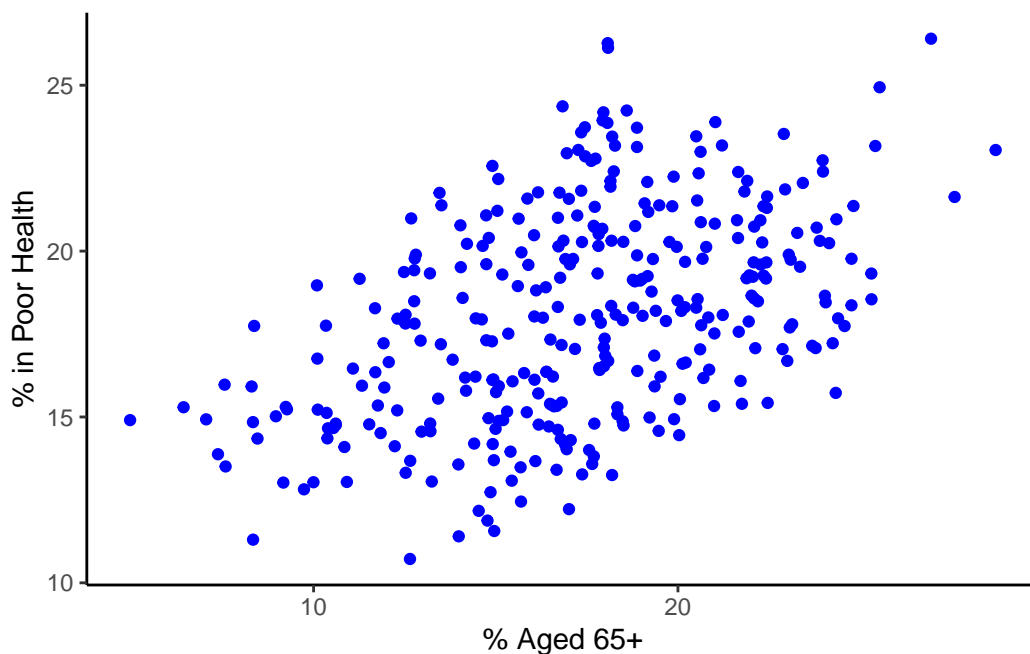
QUESTION 3

What pattern(s) do you observe comparing the spatial distribution of both variables?

4.4.3 Visualise associations with a scatter plot

As we've seen in previous weeks, another way of visualising the association between two continuous variables is through a scatter plot:

```
census_map %>%  
  ggplot() +  
  geom_point(aes(x = Age_65_plus, y = Poor_Health), colour = "blue") +  
  labs(x = "% Aged 65+", y = "% in Poor Health") +  
  theme_classic()
```



This graph reinforces the impression that older areas tend to have more people in poor health. However, the question still remains, how strong is the association/relationship between these two variables? Remember, as discussed in class, if this was a 'perfect' relationship between both variables that those dots would form a straight line. This is not the case but there still seems to be a relationship there. Discuss with your colleague next to you how strong this relationship and see if you agree.

4.5 Measuring the correlation of *continuous* variables

As covered in this week's lecture, we can quantify the amount of association between two continuous variables by calculating their correlation. The correct measure of correlation to use

(Pearsons or Spearmans) depends on the distribution of the two variables being correlated.

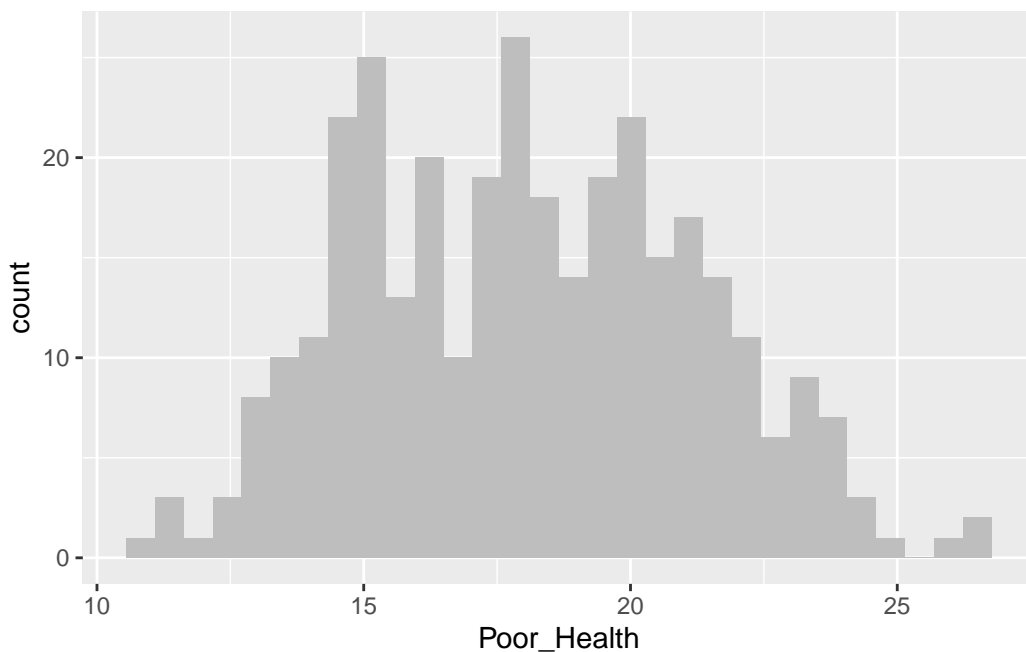
4.5.1 Check if both variables are symmetrical

If both variables are continuous and symmetrically distributed, then the correct measure of correlation to use is Pearson's product moment correlation coefficient. Alternatively, if one or both of the variables are skewed, the correct measure to use is Spearman's Rank correlation coefficient.

We can check whether or not *Poor_Health* and *Age_65_plus* are symmetrically distributed by visualising their distributions using histograms, and by calculating their skew.

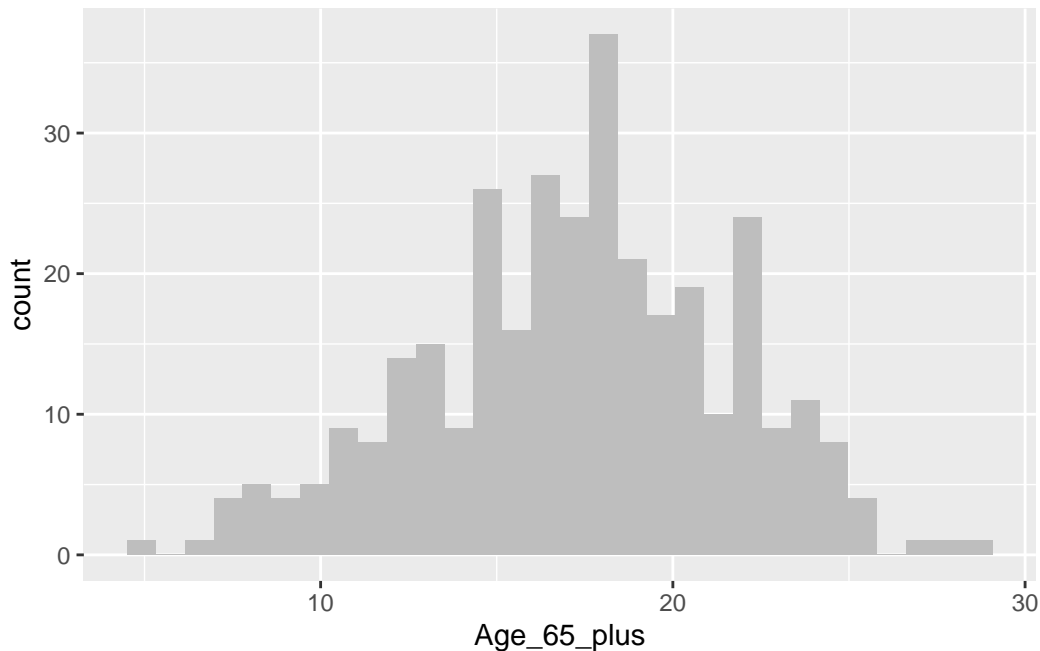
Histogram for *Poor_Health*

```
census_map %>%  
  ggplot() +  
  geom_histogram(aes(x = Poor_Health), fill = "grey", bins = 30)
```



Histogram for *Age_65_plus*

```
census_map %>%  
  ggplot() +  
  geom_histogram(aes(x = Age_65_plus), fill = "grey", bins = 30)
```



As we've seen previously we can use a boxplot to help us determine whether a variable may be skewed or not. However, sometimes this can prove challenging. In cases like this, we can compute a number that tells us how skewed a variable is. In this case, we'll measure skew in the range of -1 to +1, with values closer to 0 being less skewed, and values closer to -1 or +1 being a lot more skewed.

Compute skew for *Poor_Health*

```
x <- census_map$Poor_Health
x <- x[!is.na(x)]
mean((x - mean(x))^3) / sd(x)^3
```

```
[1] 0.1277937
```

Compute skew for *Age_65_Plus*

```
x <- census_map$Age_65_plus
x <- x[!is.na(x)]
mean((x - mean(x))^3) / sd(x)^3
```

```
[1] -0.1794138
```


From the above we can see that although neither variable is perfectly symmetrically distributed (skew = 0), they are both approximately symmetrical (skew in the range -0.5 to +0.5).

4.5.2 Calculate appropriate measure of correlation

Both *Poor_Health* and *Age_65_plus* are approximately symmetrically distributed. Therefore the appropriate measure of correlation to use is Pearson's product moment correlation coefficient. Let's calculate this correlation coefficient now between these two variables.

```
cor(census_map$Poor_Health,  
    census_map$Age_65_plus,  
    method = "pearson",  
    use = "complete.obs")
```

```
[1] 0.4631554
```

From the output above we can see that *Poor_Health* and *Age_65_plus* have a Pearson's correlation of +0.46. The size of the correlation (0.46) indicates that there is a moderate strength association between the two variables. (A correlation 0 would indicate no association, and a correlation of +1 or -1 would indicate a perfect one-to-one association.) The sign of the correlation (+) indicates that districts with more persons aged 65+ have more poor health.

If one or both variables had been skewed, then we would have needed to calculate Spearman's correlation coefficient, which can be achieved simply by changing the measure requested from "pearson" to "spearman".

4.5.2.1 Interpretation of Correlation (r) values

| Correlation (r) | Strength | Direction | Interpretation |
|-----------------|----------|-----------|-------------------------------------|
| +1.0 | Perfect | Positive | Variables increase together exactly |
| +0.7 to +0.9 | Strong | Positive | Strong increasing relationship |
| +0.4 to +0.6 | Moderate | Positive | Clear increasing trend |
| +0.1 to +0.3 | Weak | Positive | Slight increasing trend |
| 0 | None | None | No linear relationship |
| -0.1 to -0.3 | Weak | Negative | Slight decreasing trend |
| -0.4 to -0.6 | Moderate | Negative | Clear decreasing trend |
| -0.7 to -0.9 | Strong | Negative | Strong decreasing relationship |
| -1.0 | Perfect | Negative | Variables decrease together exactly |

QUESTION 5

Which statements are correct?

1. Pearson's correlation is _____.
 - a rank-based correlation
 - a linear correlation
 - a measure of spatial autocorrelation
 - a robust estimator of median association
 2. Spearman's correlation is _____.
 - a rank-based correlation
 - a linear correlation
 - a measure of spatial autocorrelation
 - a robust estimator of median association
 3. The correlation between Age_65_plus and Poor_Health is +0.46. Does this mean ageing causes poor health? Explain.
-

4.5.3 Visualising the association between categorical variables

To work with categorical variables we will be using the SAR dataset. However, before you can use this data, you must pass the SAR Data Access test on canvas located in Week 4. this should take no more than a minute or so. After passing the test, you should be able to download the file *England and Wales SAR.csv* from the Week 4 folder on CANVAS and save it to the same folder as your Week4.R script on your M drive.

Next we need to load the data,

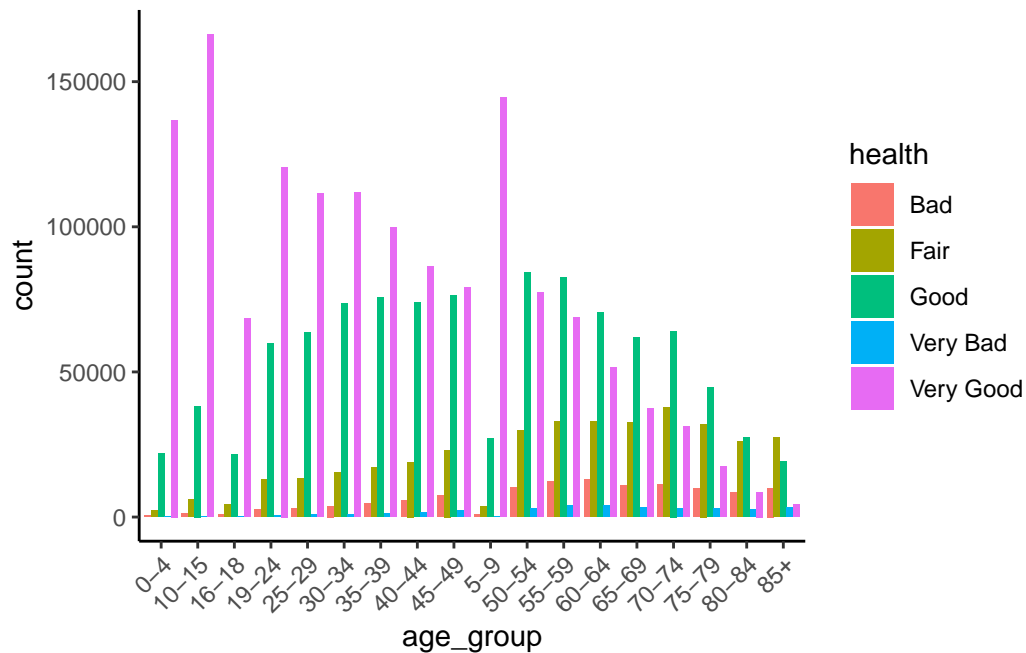
```
SAR <- read.csv("England_Wales_SAR.csv",stringsAsFactors = TRUE)
```

The *SAR* dataset contains information on an anonymised 5% sample of respondents to the 2021 Census. Anonymisation strategies range from the obvious (name and address removed) to the less obvious (e.g. the removal of all geography below district-level; the regrouping of age from single year of age to five-year age groups etc.).

The dataset contains a LOT of respondents!

Let's visualise the relationship between *health* and one of the variables that we think helps explain variations in *health*, such as *age_group*. This may take up to 30 seconds or so.

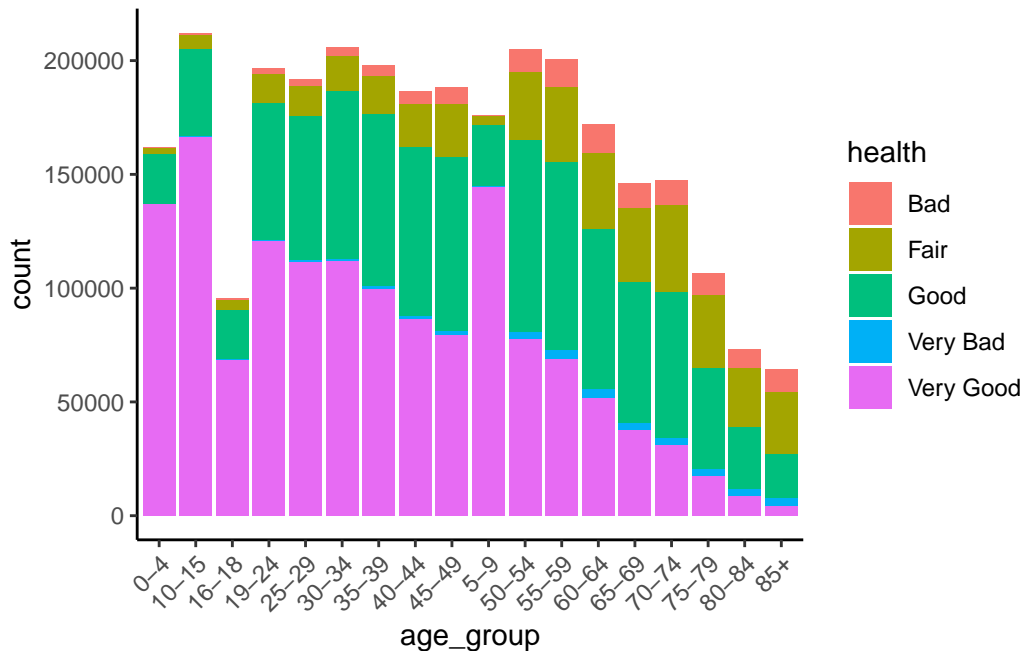
```
SAR %>%
  ggplot() +
  geom_bar(aes(x = age_group, fill = health),
    position = "dodge") +
  theme_classic() +
  theme(
    axis.text.x = element_text(angle = 45, hjust = 1)
  )
```



The clear story told by this graph is that the number and proportion of people in ‘Very Good’ health declines with age.

We can also visualise this as a stacked barchart.

```
SAR %>%
  ggplot() +
  geom_bar(aes(x = age_group, fill = health)) +
  theme_classic() +
  theme(
    axis.text.x = element_text(angle = 45, hjust = 1)
  )
```



4.6 Measuring the correlation of categorical variables

Bar graphs can give us some clues, but to be sure we need to quantify the amount of association each pair of categorical variables by calculating their correlations.

When measuring the strength of association between two categorical variables the correct measure of correlation to use is Spearman's correlation coefficient or Cramer's V. Which specific one depends on the nature of the two variables being correlated.

To select the correct measure of correlation for a given pair of categorical variables, you first need to decide whether each variable is nominal or ordinal.

```
SAR %>%
  pull(health) %>%
  levels()
```

```
[1] "Bad"      "Fair"     "Good"     "Very Bad" "Very Good"
```

```
SAR %>%
  pull(marital_status) %>%
  levels()
```

```
[1] "Divorced/partnership dissolved"  "Legally partnered, but separated"
[3] "Married"                        "Never legally partnered"
[5] "Registered Civil Partnership"    "Widowed"
```

Since both *health* and *marital_status* are not ordinal, it is clear that Cramer's V should be used.

```
tab <- table(SAR$health, SAR$marital_status)

chisq <- chisq.test(tab)

n <- sum(tab)
k <- min(nrow(tab) - 1, ncol(tab) - 1)

cramerv <- sqrt(chisq$statistic / (n * k))
cat("Cramér's V =", cramerv, "\n")
```

Cramér's V = 0.1231071

Remember that unlike Pearson's product moment coefficient and Spearman's rank, which have values in the range of -1 to +1, Cramer's V have values between 0 and 1, with 0 representing no association and 1 being a very strong association.

QUESTION 6

Which of the following correctly describes the strength of the above association?

- None
- Weak
- Moderate
- Strong

4.7 Area-level data analysis traps - MAUP

Recall that the **Modifiable Areal Unit Problem (MAUP)** is a source of statistical bias in spatial analysis whereby results (e.g., correlations, regression coefficients, summary statistics) change when the spatial units of aggregation are modified, either by altering their **scale (size)** or their **zoning (boundary configuration)**.

Let's check this relationship for *Poor_Health* and *Age_65_plus* at the Ward, Counties, District, and Region levels

First we'll need to load the Wards, Counties, and Region data. No need to load the District data because we've already done this before.

```
# Load regions file
regions <- read.csv("Census_2021_Regions.csv", stringsAsFactors = TRUE)
# Load wards file
wards <- read.csv("Census_2021_Wards.csv", stringsAsFactors = TRUE)
# Load counties file
counties <- read.csv("Census_2021_Counties.csv", stringsAsFactors = TRUE)
```

We'll use Spearman's rank correlation to quantify the association between *Poor_Health* and *Age_65_plus*

Regions.

```
cor(regions$Poor_Health,
     regions$Age_65_plus,
     method = "spearman",
     use = "complete.obs")
```

```
[1] 0.3212121
```

Counties

```
cor(regions$Poor_Health,
     regions$Age_65_plus,
     method = "spearman",
     use = "complete.obs")
```

```
[1] 0.3212121
```

Districts

```
cor(districts$Poor_Health,  
    districts$Age_65_plus,  
    method = "spearman",  
    use = "complete.obs")
```

```
[1] 0.4622558
```

Wards

```
cor(wards$Poor_Health,  
    wards$Age_65_plus,  
    method = "spearman",  
    use = "complete.obs")
```

```
[1] 0.2574719
```

In summary, the results above show that the correlation of *Poor_Health* with *Age_65_plus* vary by scale as follows:

| | Scale | Spearman |
|---|----------|----------|
| 1 | Region | 0.32 |
| 2 | County | 0.51 |
| 3 | District | 0.46 |
| 4 | Ward | 0.26 |

The correlation is strongest at county level, and decreases with the size of the spatial unit from county to ward. This suggests that at smaller scales, such as wards, there is some process that weakens the relationship between age and health. This process is likely to be some kind of social sorting (e.g. spatial segregation of the rich and poor) which results in some wards containing an unusually high share of healthy or unhealthy elderly. In contrast, there is less social sorting between counties, as each tends to contain a wide variety of persons. The weaker correlation at regional level is probably due in part to how few regions there are, and in part to the stark health and economic divide between regions in the South East of England and the rest of the country.

QUESTION 7

Can you think of any other factors that may help explain the differences in correlation values with scale?

4.8 Formative tasks

Task 1

Visualise the association between *Poor_Health* and *Density* in the **Wards** data for Liverpool

Task 2

Compute the right correlation measure to association between *Poor_Health* and *Density* in Liverpool. Interpret and discuss what the results mean.

Task 3

Using the above correlation metric, determine whether MAUP is an issues using all records the data for Wards, Counties, Districts, and Regions for the same variables.

Task 4

Compute the skewness between *Poor_Health* for Districts for all the data and for Liverpool only. Discuss what could be leading these results

Task 5

Calculate the interquartile range (IQR) of *Poor_Health* at the Region, District, and Ward levels and discuss the results as they relate to MAUP.

Remember:

- IQR measures the spread of the middle 50% of values.
- Larger IQR → greater inequality or heterogeneity.
- Smaller IQR → greater smoothing through aggregation