# Introduction to R

Department of Sociology, University of Oxford

Casey Breen

2023-10-05

# Welcome to "Intro to R"

- Two sessions:
  - Thursday, 1pm - 4pm
  - Friday, 9:30am - 12:30pm
- Course materials available from:
  - www.github.com/caseybreen/intro_r

# Course goals

- Why R is a powerful tool for social science research

- Install R and `RStudio`

- Introduction to R syntax and data types

- Basic understanding of data manipulation + visualization

# Course agenda

- **Session 1**
  - Introduction + installing `R` and `RStudio`
  - Overview of `RStudio` interface + `R` scripts, notebooks, quarto
  - Basic syntax and data types
  - Data import and export
- **Session 2**
  - Data manipulation (`dplyr`)
  - Data visualization (`ggplot2`)
  - Best practices: coding style, commenting, and documentation
  - Resources for self-teaching

# R and RStudio

- R is a statistical programming language
    - Download: https://cloud.r-project.org

- RStudio is an integrated development environment (IDE) for R programming
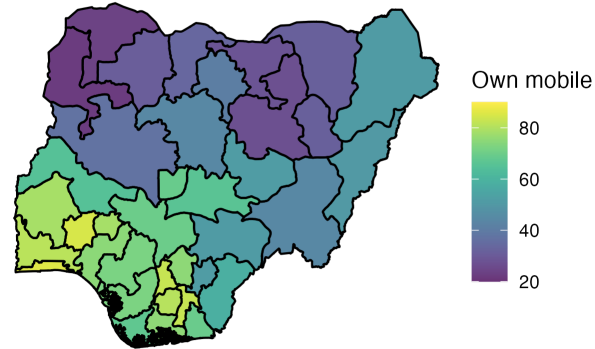    - Download: http://www.rstudio.com/download

# Why R?

- Free, open source — great for reproducibility and open science

- Powerful language for data manipulation, statistical analysis

- Publication-ready data visualizations

- Well supported, excellent community
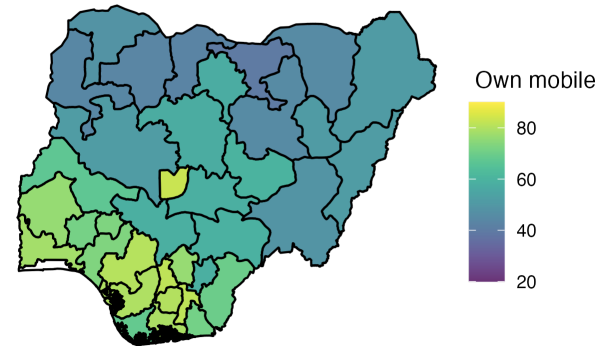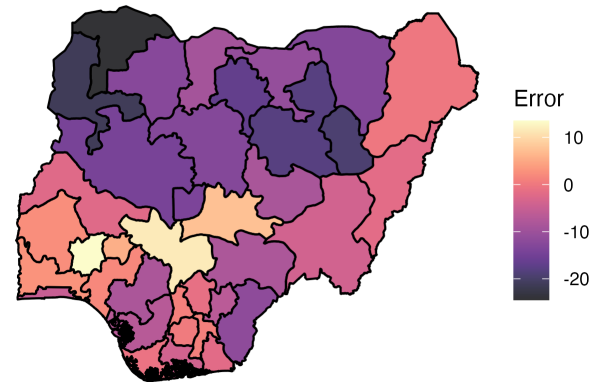
# Data visualization

# Easy to simulate + plot data

```r
1  # Generate random data for x
2  x <- rnorm(n = 10000)
3  y <- 0.8 * x + rnorm(10000, 0, sqrt(1 - 0.8^2))
4  # Create data.frame
5  data_df <- data.frame(x = x, y = y)
6  # Generate df
7  data_df %>%
8    ggplot(aes(x = x, y = y)) +
9    geom_point(alpha = 0.1) +
10   theme_classic()
```

# RStudio Panes



Scripts and Files

Enviroment

Console

Folder Tree, Viewer, Help Window, etc.

# R Scripts, R-Notebooks

- **Scripts:**
  - Just code
  - Ideal for simple tasks (and multi-script pipelines)
- **Notebooks** (Quarto, R Notebook):
  - Integrated code, text, and outputs (great for documentation!)
  - Interactive
  - We will focus on notebooks

# Quarto Document

- Notebook-Style Layout: Supports interactive code and text chunks.
    - Code Chunks: Segments for code execution
    - Text Chunks: Annotations or explanations in Markdown format.

- Inline Output: Figures and code output display directly below the corresponding code chunk

# Live Coding Session 1: Creating new Quarto file

- I'll demo first; please pay attention

# You turn

- Please create a quarto document

- You can use this document for the rest of this session

- Add a new code chunk

    - click + point: `Insert -> Executable cell -> R`

    - macOS: `Cmd` + `Option` + `I`

    - Windows/Linux: `Ctrl` + `Alt` + `I`

# Objects

- Everything in R is an object

  - **Vectors**: Ordered collection of same type.

  - **Data Frames**: Table of columns and rows.

  - **Function**: Reusable code block.

  - **List**: Ordered collection of objects.

```r
1   ## Objects in R
2
3   ## Numeric like `1`, `2.5`
4   x <- 2.5
5
6   ## Character: Text strings like `"hello"`
7   y <- "hello"
8
9   ## Boolean: `TRUE`, `FALSE`
10  z <- TRUE
11
12  ## Vectors
13  vec1 <- c(1, 2, 3)
14  vec2 <- c("a", "b", "c")
15
16  ## data.frames
17  df <- data.frame(vec1, vec2)
```

# Functions

- Built-in "base" functions

```r
1  ## Functions in R
2  result_sqrt <- sqrt(25)
3  result_sqrt
```

```
[1] 5
```

- Custom, user-defined functions

```r
1  # User-Defined Functions: Custom functions
2  my_function <- function(a, b) {
3    return(a^2 + b)
4  }
5
6  my_function(2, 3)
```

```
[1] 7
```

- Functions from packages

# Installing packages

- Packages: pre-built code and functions.

- Packages are generally installed from the Comprehensive R Archive Network (CRAN)

**Install:** new packages

```
1  install.packages("tidyverse")
```

**Library**: load installed packages

```
1  library(tidyverse)
```

**Installing a package**
```
install.packages('my.package')
```



**Loading a package**
```
library('mypackage')
```



YaRrr! The Pirates Guide to R. Nathaniel D. Phillips, 2018.

# Running code

- Run code in a quarto document (or script, or R notebook)

  - Exception: install packages, quick checks in console

- To run a single line of code

  - Cursor over line, **`Ctrl + Enter`** (Windows/Linux) or **`Cmd + Enter`** (Mac).

- To run a full code chunk (or script)

  - **`Ctrl + Shift + Enter`** (Windows/Linux) or **`Cmd + Shift + Enter`** (Mac).

# Your Turn

- Create a new code cell in `quarto` document

  - (`insert -> executable cell -> R`)

- Run each line one at a time

- Run the full code chunk

```
1  x <- 1
2
3  y <- (x + 1)^3
4
5  cat("Thank you for attending R session number", x, "!")
```
Thank you for attending R session number 1 !

# Break

10 minute break

# Basic Syntax: assignment

- Use **<-** or **=** for assignment

  - **<-** is preferred and advised for readability

- Formally, assignment means "assign the result of the operation on the right to object on the left"

```
1  ## Add comments
2
3  x <- 7 # assigns 7 to x
4
5  ## Quesiton: what does this do?
6  y <- x <- 25
```

# Basic Syntax: comments

- Use **#** to start a single-line comment

- Include lots of comments when you're writing code

```
1  ## Add comments
2
3  x <- 7 # assigns 1 to x
4  x <- 12
```

# Basic syntax: operators

```
1  ## R as a calculator (# adds a comment)
2  3 * 3
```

[1] 9

```
1  ## Division
2  12/4
```

[1] 3

```
1  ## Subtraction
2  100-12
```

[1] 88

```
1  ## Exponents (10^2)
2  10 ** 2
```

[1] 100

# Basic syntax: comparisons

| Operator | Symbol |
|---|---|
| AND | & |
| OR | \| |
| NOT | ! |
| Equal | == |
| Not Equal | != |
| Greater/Less Than | > or < |
| Greater/Less Than or Equal | >= or <= |
| Element-wise In | %in% |

```
1  ## Logical operators
2
3  10 == 10
```
```
[1] TRUE
```
```
1  9 == 10
```
```
[1] FALSE
```
```
1  9 < 10
```
```
[1] TRUE
```
```
1  "apple" %in% c("bananas", "oranges")
```
```
[1] FALSE
```
```
1  "apple" %in% "bananas" | "apple" %in% "apple"
```
```
[1] TRUE
```
```
1  "apple" %in% "bananas" & "apple" %in% "apple"
```
```
[1] FALSE
```

# Data structures

- There are lots of data structures; we'll focus on `vectors` and `data frames`.

    - `Vectors`: One-dimensional arrays that hold elements of a single data type (e.g., all numeric or all character).

    - `Data Frames`: Two-dimensional tables where each column can have a different data type; essentially a list of vectors of equal length.

# **Vectors** and **data frames**

- **Vector** example

```
1  ## Vector Example
2  vec_example <- c(1, 2, 3, 4, 5)
3
4  print(vec_example)
```

```
[1] 1 2 3 4 5
```

- **Data frame** example

```
1  # Data.frame example
2  example_df <- data.frame(
3    ID = c(1, 2, 3, 4),
4    Name = c("Alice", "Bob", "Charlie", "David"),
5    Age = c(25, 30, 35, 40),
6    Score = c(90, 85, 88, 76)
7  )
```

| ID | Name | Age | Score |
|----|------|-----|-------|
| 1 | Alice | 25 | 90 |
| 2 | Bob | 30 | 85 |
| 3 | Charlie | 35 | 88 |

# Data types

- Each `vector` or `data frame` column can only contain one data type:

  - `Numeric`: Used for numerical values like integers or decimals.

  - `Character`: Holds text and alphanumeric characters.

  - `Logical`: Represents binary values - TRUE or FALSE.

```r
1  ## generate vectors
2  vec <- c(1, 2, 3)
3  vec1 <- c("a", "b", "c")
4
5  ## check type
6  class(vec1)
```

```
[1] "character"
```

```r
1  class(vec2)
```

```
[1] "character"
```

# Generating Sequences in R

```r
1  ## Basic
2  c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)
```

```
[1]  1  2  3  4  5  6  7  8  9 10
```

- Colon operator (**:**), creates sequences with increments of 1

```r
1  c(1:10)
```

```
[1]  1  2  3  4  5  6  7  8  9 10
```

- **seq()** Function: More flexible and allows you to specify the **start**, **end**, and **by** parameters.

```r
1  ## seq 1-10, by = 2
2  seq(1, 10, by = 2)
```

```
[1] 1 3 5 7 9
```

# Basic functions

- Function: Input arguments, performs operations on them, and returns a result

- For each of the below functions, what are the:

  - Input arguments?

  - Operations performed?

  - Results?

```
1  ## generate random draws from a standard normal distribution
2  random_draws <- rnorm(n = 5,
3         mean = 0,
4         sd = 1)
5
6  ## find the mean
7  mean(random_draws)
```
```
[1] -0.0343441
```
```
1  ## find the median
2  median(random_draws)
```
```
[1] -0.1551556
```
```
1  ## find the standard deviation
2  sd(random_draws)
```
```
[1] 0.6595093
```

# In-class exercise 1

1. Assign `x` and `y` to take values 3 and 4.

2. Create a new variable `z` as the product of variables `x` and `y`.

3. Write code to calculate the square of 3. Assign this to a variable `three_squared`.

4. Write a logical expression on whether `x` is greater than 10. When might you need to filter data based on a condition?

5. Write a logical expression testing whether `x` is *not* greater than 10.

# Exercise 1 solutions

```r
1  # Assign x and y to take values 3 and 4
2  x <- 3
3  y <- 4
4
5  # Create a new variable z as the product of variables x and y
6  z <- x * y
7
8  # Write code to calculate the square of 3 and assign it to a variable three_squared
9  3^2
```

```
[1] 9
```

```r
1  # Write a logical expression to check if x is greater than 10
2  x > 10
```

```
[1] FALSE
```

```r
1  # Write a logical expression to check if x is not greater than 10
2  x <= 10
```

```
[1] TRUE
```

# In-class exercise 2

1. Generate vectors containing the numbers 100, 101, 102, 103, 104, and 105 using 3 different methods (e.g., `c()`, `seq()`, `:`). In what scenarios might each method be most convenient?

2. Generate a sequences of all **even** numbers between 0 and 100.

3. Create a descending sequence of numbers from 100 to 1, and assign it to a variable.

# Exercise 2 solutions

```r
1  # Generate a vector using c() method
2  vector_c <- c(100, 101, 102, 103, 104, 105)
3  # Use when numbers are not in a simple sequence or pattern
4
5  # Generate a vector using seq() method
6  vector_seq <- seq(100, 105, by = 1)
7  # Use when numbers follow a pattern but not necessarily just increment by 1
8
9  # Generate a vector using : operator
10 vector_colon <- c(100:105)
11 # Use when numbers increment by 1
12
13 # Generate a sequence of all even numbers between 0 and 100
14 even_seq <- seq(0, 100, by = 2)
15
16 # Create a descending sequence of numbers from 100 to 1
17 desc_seq <- seq(100, 1, by = -1)
```

# In-class exercise 3

1. Generate a sample of 100 observations drawn from a normal distribution with a mean of 10 and a standard deviation of 2. How is this type of random sampling useful in statistical analysis?

2. Calculate the mean of this generated sample. How does this sample mean relate to the population mean of the distribution?

3. Calculate the difference between the sample mean and the population mean. Why the discrepancy?

4. Repeat steps 1--3 with a sample of 10,000. Did the difference between the sample mean and the population mean decrease? Will this always be the case?

# Exercise 3 solutions

```r
1  # Generate a sample of 1,000 draws from a normal distribution with mean = 10 and sd = 2
2  sample_data_100 <- rnorm(100,
3                           mean = 10,
4                           sd = 2)
5
6  # Calculate the mean of this sample
7  sample_mean_100 <- mean(sample_data_100)
8
9  # Calculate the difference between the mean of the sample and the expected value of the mean
10 difference_100 <- abs(sample_mean_100 - 10)
11
12 difference_100
```

```
[1] 0.2247292
```

```r
1  # Calculate the Z-score for the sample mean
2  sample_data_10000 <- rnorm(10000,
3                             mean = 10,
4                             sd = 2)
5
6  # Calculate the mean of this sample
7  sample_data_10000 <- mean(sample_data_10000)
8
9  # Calculate the difference between the mean of the sample and the expected value of the mean
10 sample_data_10000 <- abs(sample_data_10000 - 10)
11
12 sample_data_10000
```

```
[1] 0.01523271
```

# Break

- 10 minutes
- Tea + cake

# Indexing vectors

- Basic indexing, specify position

```
1  vec <- c(1, 2, 3, 4, 5)
2  first_element <- vec[1]
3  third_element <- vec[3]
```

- Conditional indexing, specify position

```
1  vec <- seq(5, 33, by = 2)
2  vec[vec > 25]
```

```
[1] 27 29 31 33
```

# Working with **data frames**

- **Data frames** are the most common and versatile data structure in R

- **Data frames** are structured as rows (observations) and columns (variables)

```r
1  test_scores <- data.frame(
2    id = c(1, 2, 3, 4, 5),
3    name = c("Alice", "Bob", "Carol", "Dave", "Emily"),
4    age = c(25, 30, 22, 28, 24),
5    gender = c("F", "M", "F", "M", "F"),
6    score = c(90, 85, 88, 92, 89)
7  )
8
9  knitr::kable(test_scores)
```

| id | name | age | gender | score |
|----|------|-----|--------|-------|
| 1 | Alice | 25 | F | 90 |
| 2 | Bob | 30 | M | 85 |
| 3 | Carol | 22 | F | 88 |
| 4 | Dave | 28 | M | 92 |
| 5 | Emily | 24 | F | 89 |

# Working with **data frames**

- `head()` - looks at top rows of the `data frame`

- `$` operator - access a column as a `vector`

```
1  ## print first two rows  first row
2  head(test_scores, 2)
```

```
   id  name age gender score
1  1  Alice   25      F    90
2  2   Bob   30      M    85
```

```
1  ## access name column
2  test_scores$name
```

```
[1] "Alice" "Bob"    "Carol" "Dave"   "Emily"
```

```
1  ## all rows, columns 1-3
2  test_scores[,1:3]
```

```
   id  name age
1  1  Alice   25
2  2   Bob   30
3  3  Carol   22
```

```
4  4  Dave  28
5  5 Emily  24
```

```r
1  ## all columns, rows 4-5
2  test_scores[4:5,]
```

```
   id   name age gender score
4   4   Dave  28      M     92
5   5  Emily  24      F     89
```

# Subsetting `data frames`

- **Methods:**
  - **`$`**: Single column by name.
  - **`df[i, j]`**: Row **i** and column **j**.
  - **`df[i:j, k:l]`**: Rows **i** to **j** and columns **k** to **l**.
- **Conditional Subsetting:** `df[df$age > 25, ]`.

# Quiz

Which rows and will this return?

```
1  test_scores[1:3,]
```

- Which rows and which columns will this return?

```
1  test_scores[test_scores$score >= 90, ]
```

# Answers

```
1  test_scores[test_scores$score >= 90, ]
```

```
   id   name age gender score
1  1 Alice  25      F     90
4  4  Dave  28      M     92
```

```
1  test_scores[test_scores$score >= 90, ]
```

```
   id   name age gender score
1  1 Alice  25      F     90
4  4  Dave  28      M     92
```

# Explore `data frame` characteristics

## Check number of rows

```r
1  ## check number of rows (observations)
2  nrow(test_scores)
```

```
[1] 5
```

## Check number of columns

```r
1  ## check number of columns (variables)
2  ncol(test_scores)
```

```
[1] 5
```

## Check column names

```r
1  names(test_scores)
```

```
[1] "id"      "name"    "age"     "gender" "score"
```

# Reading in data

**Common Formats**

- CSV, Excel, TXT

**Key Functions**

- **`read.csv()`**: Read CSV files

    - Faster alternatives: `read_csv` from `tidyverse` and `fread()` from `data.table`

- **`read.table()`**: Read text files

- **`readxl::read_excel()`**: Read Excel files

```
1  ## read in CSV file
2  df <- read.csv("/path/to/your/data.csv")
3  df <- read_csv("/path/to/your/data.csv") ## faster
4
5
6  ## read in stata file
7  library(haven)
8  data <- read_dta("path/to/file.dta")
```

# In-class Exercise 4

- Let's work with a real-world, social science dataset

  - CenSoc-Numident, individual-level mortality dataset

  - https://shorturl.at/gnBQS

- Please download the CenSoc-Numident Demo file (.csv) and code (pdf) bookfrom the Harvard DataVerse

# In-class Exercise 4 (cont.)

1. Install and library the `tidyverse` package.

2. Read in the dataset using `read_csv()` from the `tidyverse` package.

3. How many columns does that dataset have?

4. How many rows the dataset have?

5. What are the column names? What type of research question could we use this dataset for?

# Exercise 4 Solutions

```
 1  1. ## install packages
 2  install.packages(tidyverse)
 3
 4  ## library tidyverse
 5  library(tidyverse)
 6
 7  2. ## read in data
 8  censoc_numident <- read_csv("/path/to/censoc_numident_demo_dataset.csv")
 9
10  3. ## nrows
11  nrow(censoc_numident)
12
13  4. ## ncols
14  ncols(censoc_numident)
15
16  5. ## column names
17  names(censoc_numident)
```

# Thank you

- Session tomorrow: 9:30am – 12:30pm

- Please try to finish exercises in advance

- Questions: casey.breen@sociology.ox.ac.uk