```
course = "Improving your statistical inferences through simulation studies in R"


lesson_iteration = 1
lesson_title = "orientation + foundational concepts"


auth = "Ian Hussey"
dept = "Psychology of Digitalisation"
```

# Why am I here?

I'm a user of stats, not a statistician or mathematician.

I'm a user of code. I'm self taught, not a Computer Science graduate or trained coder.

I use simulations to teach myself, and others, about quantitative methods to use them in research.

# Why am I here?

I'm a user of stats, not a statistician or mathematician.

I'm a user of code. I'm self taught, not a Computer Science graduate or trained coder.

I use simulations to teach myself, and others, about quantitative methods to use them in research.

## Duzen

Call me Ian

Please use a name card so I know yours 🙂

## Pronunciation & Pronouns

## ? Accessibility

Please contact me if you encounter barriers that need to be overcome
*Including English!*

## Contact

Slack where possible
ian.hussey@unibe.ch

aut = "Ian Hussey";

#

#
w

dept = "Psychology of Digitalisation || Digitalisation of Psychology"

# Why are you here?

What do you want to get from this course?

What is your existing skill level?

    Programming languages

    Confidence

What career directions interest you?

If R was an animal, what animal would it be?

# Why simulate?

It gives you access to ground truth

Take no-one's word, not even $R$'s

Helps you avoid unintentional $p$-hacking

Learn how to use a method before applying it to your real data

Significant results no longer function as a stop signal for you to consider the analysis correct/complete.

# What we will cover

Data simulation from scratch, with a focus on:

    Visibility of intermediate steps and data

    Maximising code reusability

Very little math

    Often the point of simulation is to avoid math

Lots of code

    tidyverse wherever possible

$u^b$

| # | Date | Topic |
|---|------|-------|
| 1 | 19.02.2025 | Introduction + foundational concepts |
| 2 | 26.02.2025 | Writing functions |
| 3 | 05.03.2025 | General structure of a simulation |
| 4 | 12.03.2025 | Understanding $p$-values |
| 5 | 19.03.2025 | Factorial vs. one-at-at-time simulations |
| 6 | 26.03.2025 | Hidden multiplicity in ANOVA |
| 7 | 02.04.2025 | What does it mean to violate assumptions? |
| 8 | 09.04.2025 | <<Probably no class - Ian at a conference. To be confirmed.>> Otherwise: Simulating causal models |
| 9 | 16.04.2025 | The difference between significant and non-significant is not itself significant |
|   | 23.04.2025 | No class (spring break) |
| 10 | 30.05.2025 | Understanding Confidence Intervals via sequential testing |
| 11 | 07.05.2025 | Should we test our statistical assumptions? |
| 12 | 14.05.2025 | How standardized are 'standardized' effect sizes? |
| 13 | 21.05.2025 | Meta-analysis and bias |
| 14 | 28.05.2025 | The impact of careless responding on correlations |

## SCHEDULE

| # | Date | Topic |
|---|------|-------|
| 1 | 19.02.2025 | Introduction + foundational concepts |
| 2 | 26.02.2025 | Writing functions |
| 3 | 05.03.2025 | General structure of a simulation |
| 4 | 12.03.2025 | Understanding $p$-values |
| 5 | 19.03.2025 | Factorial vs. one-at-at-time simulations |
| 6 | 26.03.2025 | Hidden multiplicity in ANOVA |
| 7 | 02.04.2025 | What does it mean to violate assumptions? |
| 8 | 09.04.2025 | <<Probably no class - Ian at a conference. To be confirmed.>> Otherwise: Simulating causal models |
| 9 | 16.04.2025 | The difference between significant and non-significant is not itself significant |
|   | 23.04.2025 | No class (spring break) |
| 10 | 30.05.2025 | Understanding Confidence Intervals via sequential testing |
| 11 | 07.05.2025 | Should we test our statistical assumptions? |
| 12 | 14.05.2025 | How standardized are 'standardized' effect sizes? |
| 13 | 21.05.2025 | Meta-analysis and bias |
| 14 | 28.05.2025 | The impact of careless responding on correlations |

$u^b$

aut = "Ian Hussey";

dept = "Psychology of Digitalisation || Digitalisation of Psychology"

The content and pacing of the course will be adapted, to some degree, to students' needs and wants. There is a selection of other topics that we could cover instead of the listed topics if you prefer, including:

- Simulating individual datasets that meet the specific experimental design of your real-world study to allow you to write your analysis code before the data is collected.
- The impact of different p-hacking strategies on false positive rates.
- The impact of different data tampering strategies on false positive rates.
- Why most psychology research is statistically unfalsifiable.
- The reliability paradox: why unreliable measures can sometimes produce replicable effects
- How confounding can produce replicable but incorrect conclusions.
- Using simulation studies to understand Bayesian estimation methods and the influence of the choice of prior.
- The impact of confusing SE and SD when extracting effect sizes for meta-analysis
- The efficacy of different methods to correct for bias in meta-analysis

# Requirements & assessment

Laptop + recent version of R, RStudio, & {tidyverse}

Weekly attendance (80% minimum)

3 at-home assignments during teaching term

    Best 2 scores count towards your grade (20% each)

1 larger assignment to be completed by <<agreed date>> (60%)

    Choose, design, implement, and report a simulation study

    Scope to be determined in class

    Start early! Ask questions!

Assignments in English (preferably) or German (if necessary)

# Requirements & assessment

All assessments must be licensed CC BY 4.0

i.e., can be used or modified with attribution

# What is difficult about this course

This course does not require you to be expert in R

But it does require that you want to *become* expert in R

You will learn about coding concepts and statistical concepts *at the same time*

Like any spoken language, 'speaking' is harder than 'understanding'
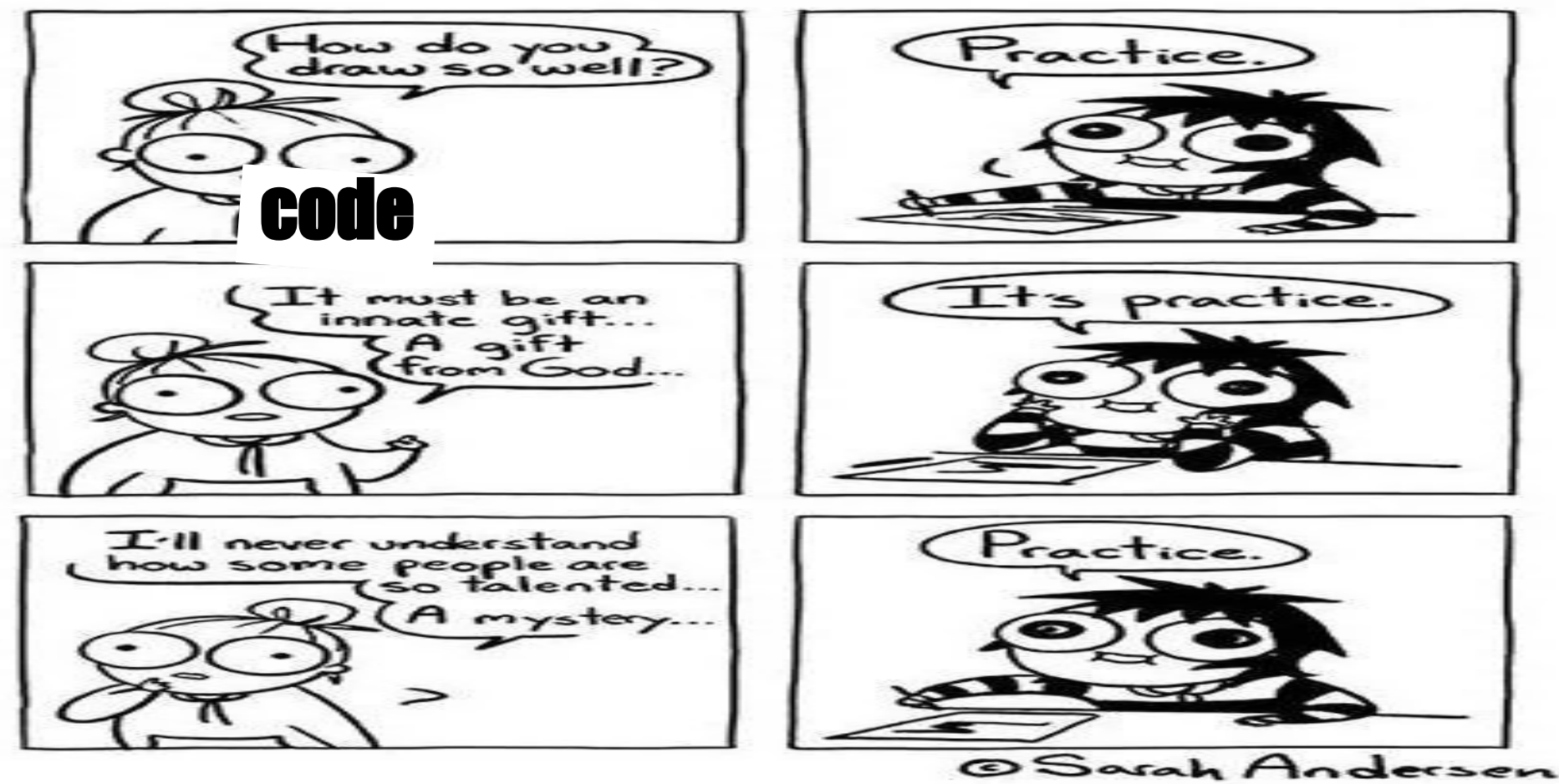
    You have to practice writing code from nothing

    Don't just read and run my code

# How to succeed in this course

Practice at home

Ask questions

Use AI (chatGPT, Gemini, codepilot, etc) the *right amount*

# What is a Monte Carlo simulation?

There is no consensus on how Monte Carlo should be defined!

Monte Carlo methods for quantitative (social) science methods research

This course

Monte Carlo methods as part of data analysis (e.g., MCMC in Bayesian data analysis)

Monte Carlo methods for the solution of general numerical problems (e.g., Monte Carlo integration)

Not this course

# **Core components** of a simulation

1. Generate pseudo-random data set with known properties
2. Analyse data with a statistical method
3. Repeat 1 & 2 many times ('iterations')
4. Summarize results across iterations
5. Make it an experiment

      Systematically vary parameters in Step 1 (between factor)

      Compare different ways to do Step 2 (within factor)

# Simulations to increase understanding

What is the distribution of *p* values under the null hypothesis?

aut = "Ian
Hussey";

dept = "Psychology of Digitalisation || Digitalisation of Psychology"

# Simulations to increase understanding

## What is the distribution of *p* values under the null hypothesis?

```r
res ← replicate(10000, t.test(rnorm(n = 50, m = 0, sd = 1), rnorm(n = 50, m = 0, sd = 1))$p.value)

res ▷ hist()
```

# Simulations to increase understanding

What is the distribution of *p* values under the null hypothesis?
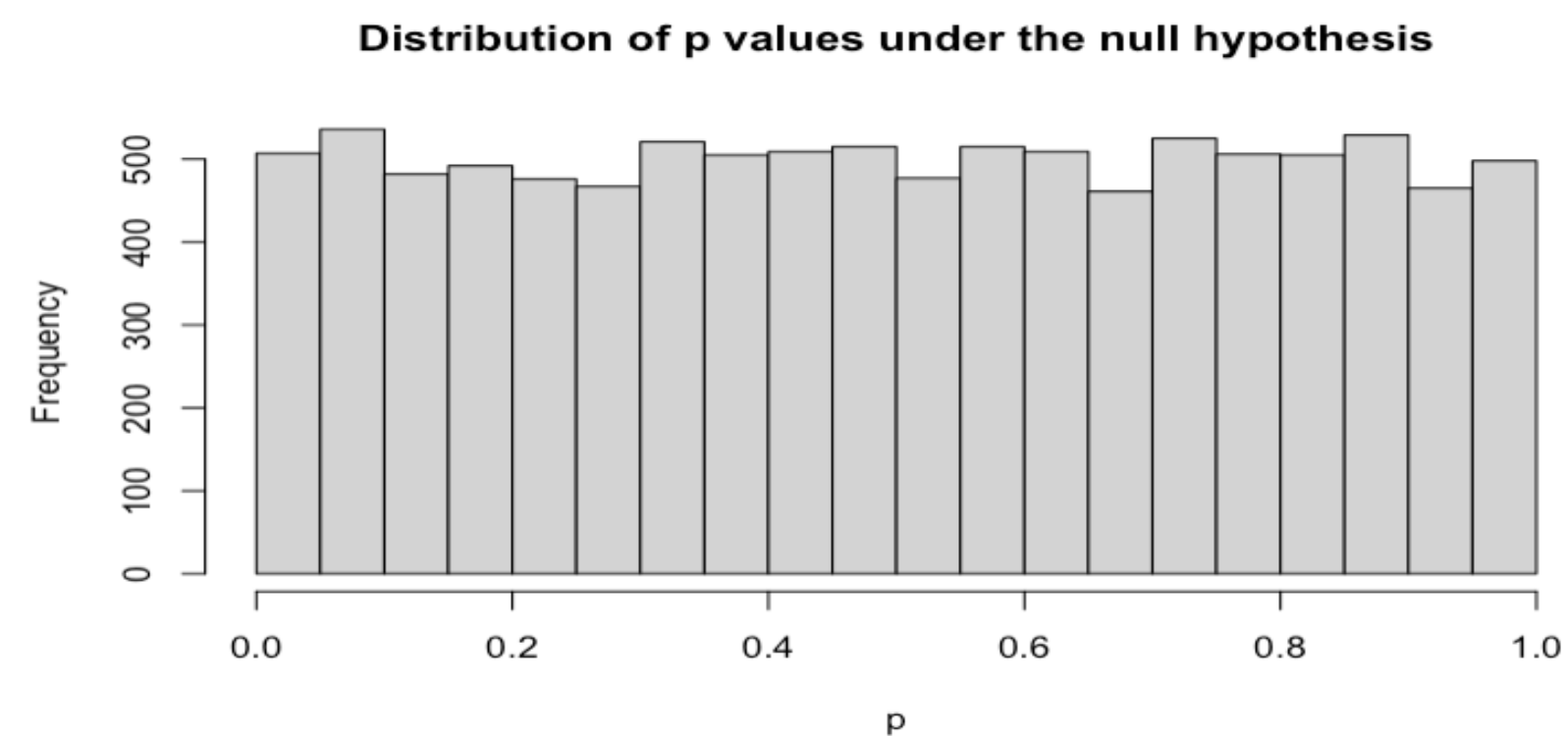
do it many
times

analyze

```
res ← replicate(10000, t.test(rnorm(n = 50, m = 0, sd = 1), rnorm(n = 50, m = 0, sd = 1))$p.value)
```

generate

generate

```
res ▷ hist()
```

summarise across iterations

aut = "Ian
Hussey";

dept = "Psychology of Digitalisation || Digitalisation of Psychology"

# Simulations to increase understanding

## What is the distribution of *p* values under the null hypothesis?

do it many times

analyze

```r
res ← replicate(10000, t.test(rnorm(n = 50, m = 0, sd = 1), rnorm(n = 50, m = 0, sd = 1))$p.value)
```

generate

generate

```r
res ▷ hist()
```

summarise across iterations

```r
1 ▾   ```{r}
2
3    replicate( # 3. repeat 1 & 2 many times ('iterations')
4      n = 10000,
5      expr = t.test( # 2. analyse data with a statistical method
6        x = rnorm(n = 50, mean = 0, sd = 1), # 1. generate pseudo-random data set with known properties
7        y = rnorm(n = 50, mean = 0, sd = 1)
8      )$p.value
9    ) ▷
10     hist(main = "Distribution of p values under the null hypothesis",
11          xlab = "p") # 4. collect and aggregate results across iterations
12
13   ```
```



Distribution of p values under the null hypothesis

# Prepare/practice

> 1_foundational_concepts__lesson.Rmd

> 2_writing_functions__lesson_and_assignment.Rmd