

## Welcome

Introduction &  
course structure

1. Introductions
2. Course motivation
3. Roadmap
4. Logistics
5. Software and computer setup
6. *Hands-on: R and RMarkdown*

McGill University is located on land which has long served as a site of meeting and exchange amongst Indigenous peoples, including the Haudenosaunee and Anishinabeg nations. McGill honours, recognizes and respects these nations as the traditional stewards of the lands and waters on which we meet today.

<https://www.mcgill.ca/fph/welcome/traditional-territory>

*see also:*

Chelsea Vowel. "Beyond Territorial Acknowledgments." Âpihtawikosisân (blog), September 23, 2016. <https://apihtawikosisan.com/2016/09/beyond-territorial-acknowledgments/>.

# Introductions

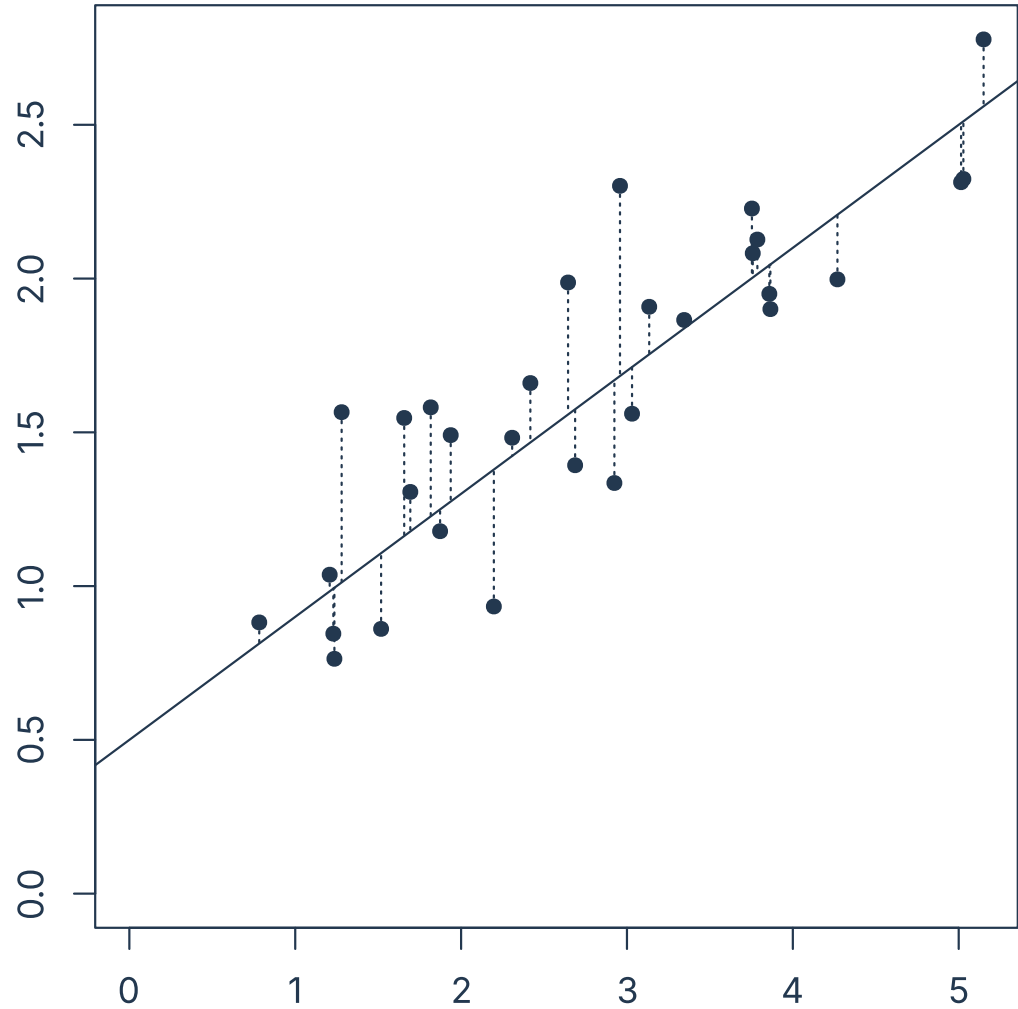


# Course motivation



Linear regression  
(OLS):

$$y_i = \alpha + \beta x_i + \varepsilon_i$$



## Linear regression (OLS):

$$y_i = \alpha + \beta x_i + \varepsilon_i$$

### **Model**

relating predictors  
to outcome

### **Assumptions**

that must be met  
for reliable  
estimation and  
interpretation

### **Estimation**

procedure to  
approximate  
unknown values

### **Language**

to talk about  
empirical "effects"

## Model

relating predictors to outcome

- ⋮ **As social scientists, the model is what we really care about**

A 'mental map' of your theoretical argument

- ⋮ **Also the fun part**

Building a tiny working model of the social world

- ⋮ **OLS (like all models) comes with *very specific* ideas about what can matter in the social world and how those things can be related**

Abbott (1988): *Transcending general linear reality*

## Estimation

procedure to approximate unknown values

- ⋮ **Predictions and measures from model and data**

- ⋮ **Technical procedures**

Important, but less *sociological*

- ⋮ **Ordinary least squares (OLS)**

- ⋮ ***But also*: maximum likelihood (ML); maximum a-posteriori (MAP); Markov chain Monte-Carlo (MCMC); ...**

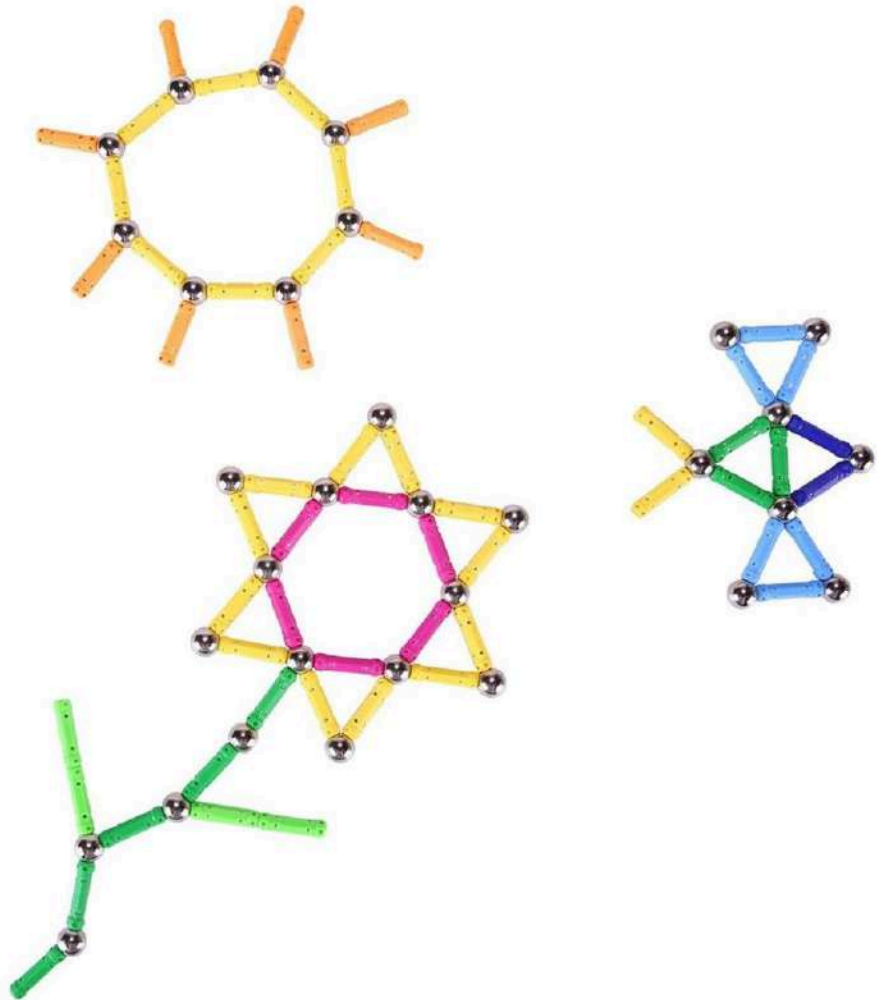
We will use the lens of *probability models* to describe all of the models in the class.

## Intuitive

- ∴ Probability distributions help to break models into components
- ∴ Probability distributions provide an intuitive language for discussing uncertainty

## Flexible

- ∴ Probability distributions describe the uncertainties in the social processes you are studying
- ∴ Simple algebra fits these distributions together to make a *model* that supports your claims





# BAYESIAN VS. FREQUENTIST STATISTICS 9

Probability models are often associated with “Bayesian” statistics, which itself is often contrasted with “Frequentist” statistics. *What do those terms mean?*

	Frequentist	Bayesian
Philosophical contrasts	<ul style="list-style-type: none"><li>∴ The <i>probability</i> of an event is the <i>proportional frequency</i> of that event across the entirety of a given ‘context’</li></ul>	<ul style="list-style-type: none"><li>∴ The <i>probability</i> of an event is a rigorous way to <i>quantify subjective uncertainty</i> about that event</li></ul>
Practical contrasts	<ul style="list-style-type: none"><li>∴ Significant limitations on <i>types of models</i> that can be used</li><li>∴ <i>Fast computation</i> of estimates for those models (OLS, ML, ...)</li><li>∴ Difficult to talk about level of <i>confidence in estimates</i></li></ul>	<ul style="list-style-type: none"><li>∴ Easy to work with a <i>wide range of models</i></li><li>∴ Estimation is <i>computationally “expensive”</i> (MCMC, Hamiltonian MC, ...)</li><li>∴ (Arguably) easy to talk about <i>confidence in estimates</i></li><li>∴ Need to specify <i>prior beliefs</i> (more on this later)</li></ul>

*In practice, these differences usually remain “under the hood.” Either approach can be used with no significant impact on reliability or credibility.*

*I strongly advocate for a pragmatic approach: use whichever framing makes the most sense for your specific model, data, resources, and audience.*

# Roadmap



## Part 1: Parametric probability models

- ⋮ Social-scientific models as random processes
- ⋮ Overview of probability distributions
- ⋮ Estimating parameters

## Part 2: Linear models and model checking

- ⋮ Re-framing linear regression as probability model
- ⋮ General model considerations (causality, overfitting)

## Part 3: Generalized linear models

- ⋮ Expanding linear models with outcome distributions and link functions
- ⋮ Binary, count, and categorical outcomes

## Part 4: Complications in data and estimation

- ⋮ Missing data and weighted observations

## Part 5: Multilevel models

- ⋮ Two-level models (nested data)
- ⋮ Covariance structures
- ⋮ Generalized multilevel models

## Part 6: Building more complex models

- ⋮ Probability models for other processes

# Logistics



## Syllabus

- ⋮ <https://soci620.netlify.app>
- ⋮ Updated regularly with links to assignments and slides and changes to the schedule

## Class periods

- ⋮ **Lecture and discussion**  
Formal discussion of topics
- ⋮ **Usually finish with demos**  
Working in R
- ⋮ **Laptop will be necessary**

## Labs

- ⋮ **Work through example code with TA**
- ⋮ **Work on assignments/projects in the same space as one another (study hall)**  
Ask questions, consult, commiserate
- ⋮ **Once per week**

## Worksheets

- ⋮ **Five worksheets over the semester**  
*Due dates on syllabus*
- ⋮ **Distributed as RMarkdown templates to complete**
- ⋮ **Everyone will evaluate *two* of their peers for each worksheet using FeedbackFruits**
- ⋮ **Turn in through MyCourses**
- ⋮ **Working together is fine (encouraged, even!), but each person needs to create their own writeup of code and expproseosition**

## Research project

- ⋮ **The main item is an original research project**
- ⋮ **Due in four parts (the four "P"s):**  
*Precis; proposal; presentation; paper*
- ⋮ **Ideally, will be part of a larger research project**  
*E.g. a draft of the methods section for a dissertation chapter?*
- ⋮ ***Meet with me early in the semester to discuss your topic ideas***



## "Generative AI"

- ∴ Language models that predict subsequent "tokens" based on previous text.
- ∴ *E.g. Microsoft Copilot (provided by McGill) OpenAI's ChatGPT, Google's Gemini, Meta's Llama, etc.*

## *The use of these tools is strongly discouraged*

- ∴ *They are bad for the world.*
- ∴ *They are bad for students.*

## Environmental impact

- ∴ Generative AI uses huge amounts electricity and water to train and to use  
"Just Five ChatGPT Queries Can Use 16oz of Water, Say Researchers."
- ∴ Generative AI contributes significantly to climate change  
"Google emissions jump nearly 50% over five years as AI use surges."



## Human exploitation

- ∴ Generative AI relies on underpaid humans to label (often harmful) content  
"What's behind the AI boom? Exploited humans."
- ∴ Generative AI is build on countless humans' uncredited, uncompensated creative work



## “Typical” text

- ∴ The technology that makes generative AI work is essentially like the predictive text on your phone, but trained on as much of the internet as corporations can get their hands on.

*One thousand Redditors (or Github projects) in a trenchcoat*

- ∴ The models are trained solely to sound *unsurprising*, not to recognize important or interesting ideas.

*“When ChatGPT summarises, it actually does nothing of the kind”*

## Writing is its own end

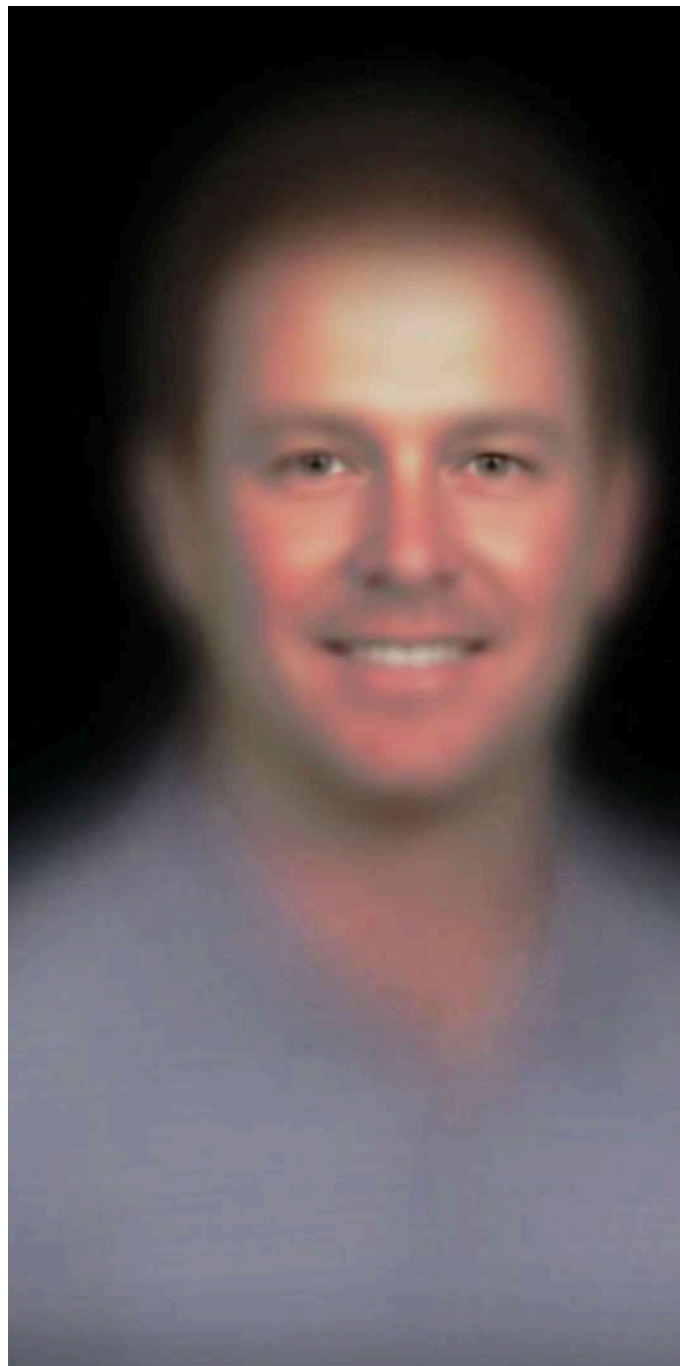
- ∴ Struggling with composition is *useful*—that is where the thinking (and learning) happens

*“We can also save time by undercooking fish, but it’s not ideal.”*

—Eryk Salvaggio

- ∴ Using AI to write your code will hinder your learning
- ∴ AI-generated code contains frequent mistakes that can be hard even for experts to spot

- ∴ *You are here to learn, and coding with AI will hinder that learning*



## Microsoft Teams

- ∴ [Available at this link](#) through browser or app
- ∴ Q&A and discussions (ask and answer!)
- ∴ Best place to contact me
- ∴ Let me know if you have trouble with access

## MyCourses

- ∴ Turning in assignments
- ∴ FeedbackFruits for peer assessment

## Readings

- ∴ Richard McElreath's *Statistical Rethinking*, (Second Edition)  
[Online access through library](#)

## The R language

- ‡ Class, labs, and worksheets will use R
- ‡ Open source (free forever)
- ‡ Vibrant ecosystem of add-on packages
- ‡ *De facto* standard for scholarly statistics

## RMarkdown

- ‡ Plain-text format to incorporate R code into documents
- ‡ Converts to Word, PDF, HTML, ...
- ‡ (*Quarto* is very similar to RMarkdown)

## RStudio (*optional*)

- ‡ A convenient interface to R and RMarkdown
- ‡ Made by Posit, the "*opinionated*" company behind `tidyverse`
- ‡ Alternatives:  
VSCoDe (VSCodium) from Microsoft;  
or any text editor and terminal

## RStudio (*or VSCode*)

User-friendly interface to the R environment and RMarkdown



## R

Statistical language and environment (the 'engine' of your analysis)



### rethink- ing

Textbook companion package



### brms

R package for Bayesian model estimation



### lme4

R package for multilevel GLM estimation



...

Other R packages (tidyverse, data.table, ggplot, ...)



## stan

General-purpose software for MCMC estimation

.10  
.01

## Installing

- ⋮ Detailed instructions to install necessary software are available at:  
<https://soci620.netlify.app/pages/software.html>

## Testing

- ⋮ A simple script to test the `rethinking` installation is at:  
[https://soci620.netlify.app/labs/lab\\_1.R](https://soci620.netlify.app/labs/lab_1.R)
- ⋮ You can download and run this, copy and paste it, or run the whole thing from directly in R:

```
source("https://soci620.netlify.app/labs/lab_1.R")
```

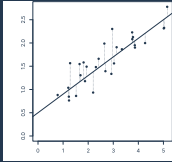
# Image credit



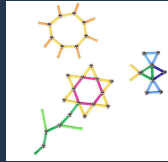
Photo by Marlis Trio Akbar on Unsplash



Photo by John Hritz on Flickr



R script to produce figure



Playmatey magnetic building blocks via WorthPoint



Coloured engraving by S.J. Neele after L. Hebert. Wellcome Collection.

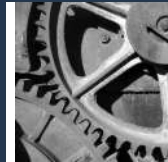


Photo by Natasha Wheatland on Flickr

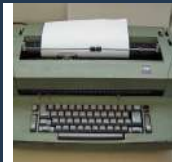
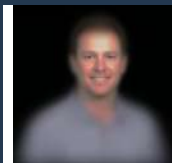


Photo by Wikimedia user Etan J. Tal



Photo by Patrick Hendry on Unsplash



Faces of 500 professional golfers, averaged by Reddit user u/osmutiar/