

TEACHING DOSSIER

Luke Hagar

Research Officer
Clinical Trials Capability
The University of Queensland

September 2025

Table of Contents

| | | |
|-----|--|----|
| 1 | Statement of Teaching Philosophy | 1 |
| 2 | Teaching Experience..... | 2 |
| 2.1 | Lecturer | 2 |
| 2.2 | Teaching Assistant | 3 |
| 2.3 | TA Workshop Facilitator | 4 |
| 2.4 | TA Coordinator | 4 |
| 2.5 | Workshops..... | 4 |
| 3 | Teaching Strategies..... | 5 |
| 3.1 | In the Classroom..... | 5 |
| 3.2 | Outside the Classroom | 6 |
| 3.3 | Assessment Methods | 6 |
| 4 | Evaluation of Teaching..... | 7 |
| 4.1 | Evaluation from Students..... | 7 |
| 4.2 | Corroboration by Faculty Members | 8 |
| 5 | Professional Development | 9 |
| 5.1 | Fundamentals of University Teaching | 9 |
| 5.2 | University Mathematics Teaching Techniques | 9 |
| 5.3 | Certificate in University Teaching..... | 9 |
| 5.4 | New Instructor Foundations Program | 9 |
| 6 | Future Goals..... | 10 |

Appendices

| | | |
|---|--|----|
| A | STAT 341 Assignment Question..... | 11 |
| B | STAT 341 Assignment Walkthrough Videos | 13 |
| C | STAT 938 Mock Consulting Session Plans..... | 14 |
| D | BIOS 612 Jupyter Notebooks | 17 |
| E | R Shiny Application for CUT Program | 19 |

1 | Statement of Teaching Philosophy

As an educator, I draw on my experience as a student to inform my teaching priorities. In doing so, I endeavour to emulate instructors that have encouraged my passion for statistics and reflect constructively on the obstacles I faced to promote a more empathetic environment for others. I personally have a strong passion for teaching statistics and embed my values into my teaching priorities. My priorities for teaching include (i) fostering a supportive environment that is conducive to learning, (ii) providing students with opportunities to actively engage with course content, and (iii) emphasizing the importance of statistical communication skills.

From course design to delivery, I strive to cultivate an environment where diverse groups of learners feel supported. I equip students for success through course delivery by introducing statistics concepts using varied techniques. Traditional math courses are often tailored to students who learn best by engaging with algebraic expressions. While mathematical theory is important, I also emphasize the use of high-quality visualizations and interactive tools to make complex statistical concepts accessible to a broader audience of students. In line with universal course design principles, I set clear expectations for both the students and myself, and I emphasize flexibility in the course structure (e.g., late submission accommodations). To help alleviate student anxiety associated with high-stakes, in-class assessments, I begin the term by reiterating the value of meaningfully engaging with our asynchronous assessments before consulting generative artificial intelligence. In doing so, students build confidence for the exams and presentations meant to evaluate their knowledge. Evidence that I create a supportive environment for my students is reflected in feedback in [Section 4.1](#) of this dossier.

Creating a supportive classroom environment is necessary to encourage in-class engagement. I incorporate active learning opportunities into my classes by dividing them into shorter segments. In between these segments, I have students participate by engaging with technology or solving short problems. When co-teaching a graduate course at McGill University, my lectures alternated between the course slides and Jupyter notebooks that I created to lead students through simulation studies and real data analyses. I prompted discussion by challenging students to reflect on the strengths and limitations of the statistical methods in these notebooks. In a large undergraduate class that I taught at the University of Waterloo, I used low-technology polling methods (e.g., thumbs up or down gestures) to facilitate answering true or false questions. Because making mistakes is a crucial part of the learning process, these questions provided a low-stakes environment for students to evaluate their understanding. As a TA Workshop Facilitator with the University of Waterloo's Centre for Teaching Excellence, I also delivered many workshops that guided other instructors to make their own lectures more interactive.

Finally, I believe it is crucial to emphasize statistical communication in my teaching. In my courses, I require that students go beyond obtaining the correct numerical solution to a statistical problem: I ask them to interpret and communicate their findings with technical and nontechnical language. By reflecting on my interdisciplinary collaborations, I have gained a deep appreciation for why strong communication skills are important. For example, the ability to explain statistical topics to stakeholders at all levels has been incredibly valuable in my work with Airbnb. When I share such experiences with my students, they seem genuinely motivated to

take this part of their education seriously. To prepare my students for interdisciplinary work, I provide them with opportunities to develop their communication skills alongside their technical ones. I design my courses with the goal of empowering students to tackle challenging problems and effectively communicate their solutions to earn buy-in from decision makers.

I aim to provide students with experiences they can draw on after my courses by supporting a diverse student population via course design and delivery, providing regular active learning opportunities, and prioritizing statistical communication. While student performance in my courses is important, my main motivation to incorporate these values into my teaching is so that students can harness the skills they learned to effectuate change in their fields upon graduation.

2 | Teaching Experience

2.1 Lecturer

I have previously been hired on fixed-term contracts to teach two courses as a lecturer. My experience as a lecturer spans two institutions, undergraduate and graduate courses, and various subfields within statistics. At the University of Waterloo, I taught STAT 341 (Computational Statistics and Data Analysis) in the Winter 2024 term. At McGill University, I co-taught BIOS 612 (Advanced Generalized Linear Models) with Dr. Shirin Golchi in the Fall 2024 term. In this subsection of my dossier, I overview my contribution to these two courses.

2.1.1 STAT 341 – Computational Statistics and Data Analysis

Computational Statistics and Data Analysis is a third-year undergraduate course at the University of Waterloo; it is the largest and most popular elective course taken by students in the Department of Statistics and Actuarial Science. The course covers data visualization techniques, optimization algorithms, inference based on sampling distributions, and prediction methods. I taught a single 125-student section of STAT 341; however, I assumed coordination and administrative responsibilities for both sections offered in the Winter 2024 term (250 total students). In this role, I supervised six teaching assistants (TAs). I also drove the creation of assessments and their marking rubrics for all components of the course: four assignments, two midterms, and one final exam. I elaborate on one assignment that I created for STAT 341 and how it aligns with my teaching priorities in [Section 3.3](#) and [Appendix A](#) of this dossier.

Moreover, I recorded walkthrough videos for all four assignments that students in both sections could optionally view to get clarification on the expectations for each assignment question. Links to these walkthrough videos are shared in [Appendix B](#). I also held two drop-in office hours per week to assist students, and I answered over 550 student questions on Piazza, our course discussion board. Feedback that my office hours and responsiveness on Piazza made my students feel supported is provided in [Section 4.1](#). As detailed on my CV, I was awarded the Scotiabank Scholarship for my administrative work in this course.

Date: Winter 2024

Audience: 125 undergraduate students (in person)

2.1.2 BIOS 612 – Advanced Generalized Linear Models

Advanced Generalized Linear Models is a popular graduate course in McGill University's Department of Epidemiology, Biostatistics and Occupational Health. This course reviews marginal and conditional regression models for correlated continuous and categorical data in both Bayesian and frequentist frameworks. I co-taught this course with Dr. Shirin Golchi. The teaching responsibilities were divided such that I lectured in the second half of the term, and both Dr. Golchi and I were involved in creating all assessments and providing feedback.

For the Fall 2024 offering of this course, Dr. Golchi and I aimed to revamp how students interacted with R code during class. In previous terms, students were provided with static R code and output in the course slides. Students in the Fall 2024 lectures actively engaged with Jupyter notebooks that contained updated code and examples. I created these notebooks for the second half of the course. Completed versions of these notebooks are shared in [Appendix D](#). Students also had a paper presentation and a final presentation in this course. Aligning with my teaching priorities, I provided positive and constructive written feedback related to communication during the paper presentations to equip the students for success with their final presentations.

Date: Fall 2024

Audience: 10 graduate students (in person)

2.2 Teaching Assistant

At the University of Waterloo, I was a teaching assistant for several undergraduate and graduate courses that span the Department of Pure Mathematics and the Department of Statistics and Actuarial Science: *Honours Calculus I/II*, *Forecasting*, *Financial Statement Analysis*, *Computational Statistics and Data Analysis*, *Experimental Design*, and *Statistical Consulting*. Across these roles, I held drop-in office hours, graded assignments and exams, and led tutorials. Here, I highlight my involvement in one course that informed my teaching priorities.

2.2.1 STAT 938 – Statistical Consulting

Statistical Consulting is an elective graduate course in the Department of Statistics and Actuarial Science. This course overviews the technical and soft skills that are useful for statistical consultants including problem-solving techniques, effective questioning strategies, the generation of high-quality graphics, time management skills, and oral and written communication proficiency.

As a TA for the course, I helped run tutorials. Several tutorials were allotted to mock consulting sessions, where each student had 10 minutes to gain hands-on experience leading a short meeting as a consultant. I created the scenarios for those sessions, drawing on my experience as a consultant with the University of Waterloo's Statistical Consulting and Survey Research Unit. I also acted as the client during these sessions. I elaborate on the scenarios created for these sessions in [Appendix C](#). This was the first term that STAT 938 students were assessed via mock consulting sessions, and the instructor Dr. Nathaniel Stevens incorporated my work into these

sessions in subsequent iterations of the course as well. I won the Department of Statistics and Actuarial Science's 2023 Best TA Award for my contributions to this course.

Date: Spring 2023

Audience: 15 graduate students (in person)

2.3 TA Workshop Facilitator

I worked part time with the University of Waterloo's Centre for Teaching Excellence (CTE) as a Teaching Assistant Workshop Facilitator (TAWF) for four terms during my PhD. Each term, I facilitated five 90-minute workshops for the CTE's **Fundamentals of University Teaching** program, sometimes with a co-facilitator. I (co)-facilitated the following workshops: *Building TA-Instructor Rapport* (x3), *Collecting and Using Feedback on Your Teaching* (x1), *Giving and Receiving Feedback* (x4), *Supporting Student Mental Health* (x2), *Interactive Lectures* (x3), *Teaching Methods* (x2), and *Teaching STEM Tutorials* (x3). Each time that I delivered a workshop, I updated it to incorporate participant feedback from past offerings or changes to the University of Waterloo's TA policies. In addition to introducing workshop content, TAWFs guide participants through multiple interactive exercises. For several of these workshops, my co-facilitator was a new TAWF, and I contributed to their orientation.

Date: January 2023 – April 2024

Audience: 8-24 graduate students (in person and online)

2.4 TA Coordinator

Due to my prior experience as a TAWF (see [Section 2.3](#)), I was personally recruited as the inaugural TA Coordinator in the Department of Statistics and Actuarial Science at the University of Waterloo. In this part-time role, I formally contributed to the effort to revamp the Department's TA training program. Under the supervision of Dr. Chelsea Uggenti, I co-developed the final practicum component of a TA training program that students must complete before instructing their first course. Students are guided through this practicum by completing a course on the University of Waterloo's learning management system, which I designed and created. Furthermore, I conducted teaching observations to assess several TAs who were leading tutorials, and I delivered an introductory training session to 30 first-time TAs in the Department of Statistics and Actuarial Science.

Date: September 2023 – December 2023

2.5 Workshops

2.5.1 University of Toronto Health Data Working Group

I was recently invited by Dr. Kuan Liu to give a [workshop](#) for the Health Data Working Group at the University of Toronto's Dalla Lana School of Public Health. The topic of the workshop was *Sample Size Determination for Bayesian Clinical Studies*, a subject that I have considered extensively in my research program. Because of her familiarity with my research through the

Statistical Society of Canada, Dr. Liu suggested that a simplified version of my work would be well suited for the workshop's audience of master's and junior PhD students.

Before the workshop, I provided students with an R Markdown file to generate sampling distributions of posterior probabilities and assess operating characteristics for simple Bayesian designs. My research is most useful for complex designs, but it was important that students be able to run simulations quickly during the workshop. Students discovered theoretical results about study design through simulation (e.g., we explored how the sampling distribution of posterior probabilities is asymptotically uniform under null hypotheses), and we discussed the implications of these results. We also constructed visualizations of how my research methods work for simpler models, with the goal of encouraging attendees to view Bayesian design as both a practically useful tool and a promising avenue for their research.

Date: April 2025

Audience: 13 graduate students (in person)

3 | Teaching Strategies

3.1 In the Classroom

As an instructor, I provide students with opportunities for in-class engagement. I believe that actively engaging with course content helps students develop the confidence to apply statistics concepts on their own – both on course assessments and later in their chosen fields. I design my lectures to maximize the impact of interactivity on student learning.

Because students have more bandwidth to engage in class when technical content is sufficiently motivated, I begin each lecture with a short review of the preceding class. I then explicitly provide the learning objectives for the current lecture. When doing so, I preview how the upcoming lecture builds on or addresses gaps in the content that was learned previously.

Once the new content is clearly motivated, I divide my lectures into roughly 15-minute segments separated by active learning activities. I used this technique to structure my BIOS 612 classes at McGill University, where the students engaged with Jupyter notebooks (see [Appendix D](#)) between these lecture segments. In larger undergraduate classes, I often divide my lecture segments by conducting comprehension checks via multiple choice or true or false questions. After my STAT 341 students at the University of Waterloo struggled with the true or false questions on our first midterm, I incorporated more of these questions into my lectures to give students opportunities to evaluate their understanding and make errors in a low-stakes environment. Once I communicated the motivation behind participation, students responded more frequently and generally performed better on the final exam's true or false questions.

While I often conduct these comprehension checks using low-technology polling techniques, technology-enhanced polling methods could facilitate more anonymous participation. I use free polling platforms, such as *PollEverywhere*, wherever possible to reduce barriers to participation, especially if these activities are not incorporated into students' grades. If students are asked to solve a problem individually, I often give them a minute to check their answer with a partner

before responding to the poll. I find that giving students the opportunity to vet their answer with a peer gives rise to higher participation rates.

I will also seek out active learning techniques for future courses that are suitable for the students' academic level and the course structure, such as jigsaw activities or case studies. I will consult the resources provided by institutional teaching support units to expand my repertoire of active learning methods.

3.2 Outside the Classroom

I firmly believe that learning extends beyond the classroom environment. Mastery of course content does not occur in the allotted lecture time. I therefore make myself available to support my students outside the classroom. I hold weekly office hours for students to ask questions. In large classes, I also create a course discussion board to streamline the process of answering questions related to course content. My teaching evaluations in [Section 4.1](#) corroborate the importance that I place on being available to my students in and outside the classroom.

However, I recognize that my courses are not the only responsibility that my students must manage. As such, I plan to provide clear and relevant practice problems for students to engage with on their own time. By doing so, motivated students will have the necessary resources to better master the content in my courses.

It is also important to consider how these practice problems are incorporated into the classroom. Practice problems arise naturally when there is not time to complete an example at the end of class. In those scenarios, I will ask students to confirm they can obtain the final answer, which I will provide, on their own. When reviewing content at the beginning of the next lecture, I may ask students to discuss their solution with a peer before quickly debriefing the problem.

3.3 Assessment Methods

When designing my courses, I incorporate varied assessments that are aligned with my course's learning outcomes. I use these assessments to monitor students' progress in mastering the skills I have set out to teach in the course. My assessments are also informed by Bloom's taxonomy since it is important for students to demonstrate their knowledge of statistics at multiple levels of cognitive complexity. I evaluate the lower levels of Bloom's taxonomy via in-class exercises or the first questions of a midterm or final exam. These in-class exercises are diagnostic assessments that allow me to monitor whether students understand key concepts and can apply them in standard contexts throughout the term. With the current capabilities of artificial intelligence, I prefer to assess recall and understanding during class time.

I use assignments to evaluate the higher levels of Bloom's taxonomy. When completing cognitively complex assignments, it is not enough for students to memorize statistical procedures: they must practice statistical thinking. I encourage students to practice statistical thinking by requiring them (i) to write their own code to analyze or simulate data in *R* and (ii) to interpret and communicate their statistical findings with technical and nontechnical language.

These practices align with recommendations for undergraduate statistics courses prescribed in the Guidelines for Assessment and Instruction in Statistics Education (GAISE) Report.

Appendix A details one such assignment question that I created for STAT 341. In the previous question, students created an *R* function that takes a sample of data and overlays a histogram and a kernel density estimate that they computed from scratch on the same plot. In this question, students used that function with a data set detailing the number of weeks that various #1 songs spent atop the Billboard Hot 100 chart. Students determined whether the histogram or density estimate better summarizes the data (most songs spend the minimum 1 week atop the chart). This question helped students recognize that the key consideration is not whether a method *can* be applied to a data set, but whether it *should* be. The question also accessibly introduced a new concept of rejection sampling with high-quality visualizations. Students used this concept to estimate the probability of a #1 song spending less than 1 week atop the chart as 30% according to the kernel density estimate.

While my assignments are more difficult than my exams, I still evaluate the higher levels of Bloom's taxonomy on a few exam questions to distinguish between students who are meeting and exceeding course expectations. For large classes, my courses culminate with a final exam; I believe this is the most objective way to evaluate students given the time constraints with marking final assessments. For more advanced courses with fewer students, I prefer to use final projects where students must both conduct simulations and present their findings in an oral presentation or written report. These projects allow students to better apply the topics they investigate in future academic or professional contexts. I am also interested in incorporating peer assessment into my future graduate courses, where I review effective feedback strategies and evaluate students on the quality of the feedback that they give each other.

4 | Evaluation of Teaching

4.1 Evaluation from Students

The end-of-course evaluations at the University of Waterloo have students respond to several general prompts about the instructor using a scale from 1 to 5 (*Strongly Disagree* = 1, *Disagree* = 2, *Neutral* = 3, *Agree* = 4, *Strongly Agree* = 5) and answer various short-answer questions. When I taught STAT 341, 29 of 120 students who completed the course (24.2%) completed the course evaluation. The most relevant general prompts and the average score assigned to me by all respondents are provided below.

- The instructor helped me understand the course concepts (4.2)
- The instructor created a supportive environment that helped me learn (4.4)
- The intended learning outcomes were identified (4.4)
- The intended learning outcomes were assessed through my graded work (4.2)

The end-of-course evaluations at McGill University have students respond to different general prompts using the same 5-point scale along with several short-answer questions. 4 of the 10 of my BIOS 612 students (40%) completed the evaluation. The most relevant general prompts and the average score assigned to me by all respondents are provided below.

- Overall, this instructor is an excellent teacher (4.3)
- Overall, I learned a great deal from this instructor (4.5)
- The instructor conducted class sessions in an organized manner (4.5)
- The instructor's explanations were clear and understandable (4.8)
- The instructor was helpful when I sought advice (4.5)

My full teaching evaluations are available upon request. The numerical summaries above are supplemented by the following selected short-answer responses.

"Going into this course I was intimidated because I heard from friends who have taken previous iterations of this course about how hard this course was and how long and difficult the assignments are. However, thanks to [Luke's] help and teaching this term, I not only learned a lot but I never felt I was struggling or overwhelmed. The care [Luke] show[s] to the students during office hours and on Piazza has been a game changer and has made this one of my favourite courses I have taken!" – Student A (STAT 341)

"The office hours with Luke and the TA help on Piazza were extremely helpful." – Student B (STAT 341)

"Luke was quick and helpful in responding to emails about course questions whether that was for logistic reasons or help-related reasons." – Student C (BIOS 612)

"Luke comes prepared, does not rush, and knows his material well which allows for interesting class discussions." – Student D (BIOS 612)

Luke "made it super clear what the learning objectives and bigger picture was with the material. It felt like everything was related so it was easy to see the motivation from topic to topic, giving me a strong intuition about why we use one method of analysis over another." – Student E (STAT 341)

These comments and numerical summaries indicate that my teaching priorities are effectively reflected in students' learning experiences. Students A, B, and C wrote about my ability to create a supportive environment. Student D mentioned interesting class discussions, which are associated with active learning. Student E discussed developing intuition for choosing between statistical methods, and articulating this intuition is a crucial part of statistical communication.

4.2 Corroboration by Faculty Members

My references can also corroborate my teaching capabilities, particularly Dr. Nathaniel Stevens and Dr. Shirin Golchi. The first time that an instructor teaches a course in Waterloo's Department of Statistics and Actuarial Science, they are paired with a mentor who has previously taught the course. This mentor shares expertise as needed. Dr. Stevens was my mentor for STAT 341 and the instructor when I TA'd STAT 938; he therefore provides informed commentary on my contributions to those two courses. Since I co-taught BIOS 612 at McGill University with Dr. Golchi, she provides an expert perspective on my work in that course.

5 | Professional Development

5.1 Fundamentals of University Teaching

In 2021, I completed the **Fundamentals of University Teaching (FUT)** program through the CTE at the University of Waterloo. In this program, I engaged with six workshops on developing rapport with instructors, supporting student mental health, implementing active learning strategies, and giving and receiving effective feedback. I later facilitated these workshops as a TAWF (see [Section 2.3](#)), which reinforced my knowledge of their content. I also delivered three microteaching sessions, each consisting of a brief interactive lesson followed by feedback from peers in other disciplines. The FUT program is interdisciplinary, and I found it helpful to broaden my perspective on teaching by learning from peers in different fields.

5.2 University Mathematics Teaching Techniques

In the Winter 2023 term, I completed the **University Mathematics Teaching Techniques** course offered by the University of Waterloo's Faculty of Mathematics. This course was run as a weekly seminar. In this course, participants learn about active learning strategies, complete three microteaching sessions, design an assessment, deliver a guest lecture, and discuss how to handle common situations in the classroom. This seminar was specific to teaching mathematics and offered complementary perspective to that gained from the FUT program.

5.3 Certificate in University Teaching

In 2023, I completed the **Certificate in University Teaching (CUT)** program offered by the University of Waterloo's CTE. The one-year program is comprised of three courses, each of which explores a particular facet of instruction in post-secondary classes. The first course consists of four modules on student learning, interactive teaching activities, assessment methods, and course design. The second course requires participants to complete a pedagogical research project and teaching dossier. The third course is comprised of two teaching observations by members of the CTE or the participant's faculty.

For my research project, I designed a one-hour workshop for graduate students in the Department of Statistics and Actuarial Science on "Active Learning Strategies with p -Values". The workshop involved an interactive activity facilitated via an *R Shiny* app I created, which aligns with the emphasis on active learning strategies in my teaching statement. The interactive activity and link to the app is included in [Appendix E](#). This app reflects my ability to develop technology for classroom engagement. The program coordinators were so impressed with my workshop materials and corresponding literature review that my project is now a resource for future students completing this course.

5.4 New Instructor Foundations Program

In 2023, I also completed the **New Instructor Foundations Program** offered by the University of Waterloo's CTE. This three-day workshop is offered to students and postdocs who are preparing to teach their first course. The workshop covers universal design principles, syllabus

creation, and best practices to connect students with campus support services, such as the University of Waterloo's AccessAbility Services and Writing and Communication Centre.

6 | Future Goals

I am excited to take part in future teaching opportunities within the field of statistics. In 2024, I was selected to give a conference talk on my experience creating engaging assessments in undergraduate statistics courses. I hope to continue to engage with such teaching events and positively represent both the field of statistics and my future employer. As I gain more teaching experience, I also hope to teach and (re)design courses on experimental design, computational statistics, and Bayesian inference.

Throughout my teaching career, I will continue to adopt new teaching techniques that promote statistical communication and student interactivity. As a statistician, I would recommend that a collaborator or client iteratively make small changes when designing an experiment and evaluate their impact; I plan to follow a similar approach with my teaching. To ensure these changes are grounded in pedagogy, I will engage with and contribute to the workshops organized by prominent statistics education groups, including the Statistical Society of Canada's Statistical Education Section and the American Statistical Association's Section on Statistics and Data Science Education. I will supplement the knowledge from these workshops by consulting the statistics education literature, of which my preferred journals are *Journal of Statistics and Data Science Education* and the "Teacher's Corner" section of *The American Statistician*. This pedagogic awareness will also inform how my teaching priorities evolve over time.

A | STAT 341 Assignment Question

This assignment question from STAT 341 is elaborated on in [Section 3.3](#) of this dossier.

QUESTION 4: R Analysis [21 points]

The *Billboard* Hot 100 is the standard chart used by the American music industry to assess the popularity of songs. The chart has been published weekly by *Billboard* Magazine since August 4, 1958. Each week, the chart ranks the 100 most popular songs based on data from the corresponding tracking period. The chart rankings are based on radio play, online streaming, and physical and digital sales in the United States.

The archive of songs that reach #1 on the chart is well documented (see e.g., [this list](#) for the year 2023). You will work with data from songs that went #1 on the *Billboard* Hot 100 during the $N = 209$ weeks between January 1, 2020 and December 31, 2023. These data are available in the `bh100.csv` file. Each row of this file corresponds to a specific week and the columns and their contents are described below.

| Column | Description |
|--------|--|
| Year | An integer between 2020 and 2023 signifying a year. |
| Week | An integer between 1 and 53 denoting the week of the year in which the <i>Billboard</i> Hot 100 chart for that row was published. |
| Title | A character string corresponding to the title of the #1 song on the chart that week. Each song is uniquely identified by its title since no <i>distinct</i> songs that went #1 between 2020 and 2023 had the same title. |
| Artist | A character string signifying the artist(s) credited on the #1 song. |

- [2 points] Using R, read in the data found in `bh100.csv` and create a data frame with a row for each song that indicates how many weeks it spent at #1 on the *Billboard* Hot 100 between January 1, 2020 and December 31, 2023. Then, use the `summary()` function on this variable for the number of weeks and output the results.
- [4 points] Let the *population* variance of the number of weeks spent at #1 be the attribute of interest so that $a(\mathcal{P}) = SD_{\mathcal{P}}(y)$. The influence of song u on $a(\mathcal{P})$ is $\Delta(a, u)$ from the course notes. Construct an influence plot of Δ vs. the observation number. Identify the song(s) with the largest influence on the population variance attribute and determine their title(s). Based on the data frame you created in part (a), describe why the song(s) have such a large influence.
- [4 points] Using the `hist_density_plot()` function you created in Question 3, construct a plot that visualizes the histogram and density estimate for the the number of weeks spent at #1 variate. The bandwidth `h` should be determined using Silverman's "rule of thumb" introduced in Question 3 and `xrange` should be `c(-1, 18)`. Be sure to label the axes and titles informatively. Given your plot, does the histogram or density estimate better summarize the data, and why?

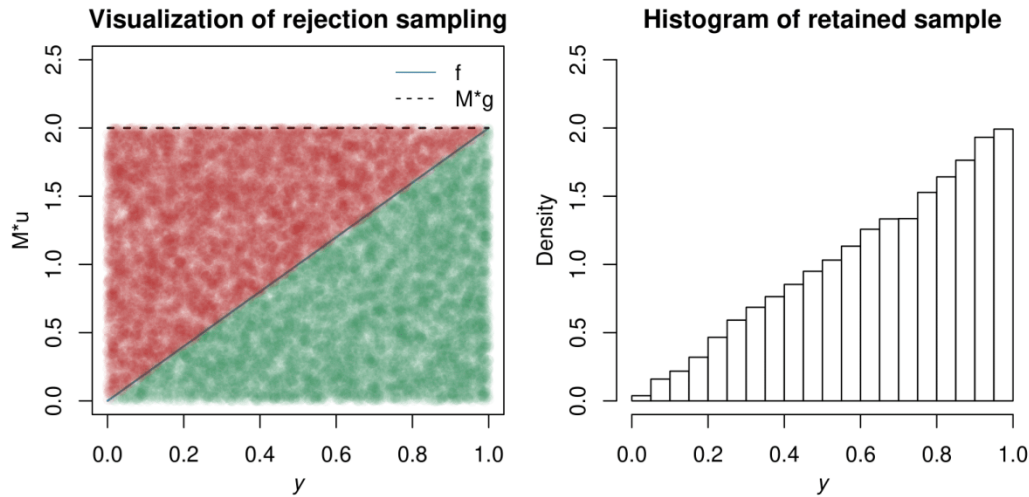
In the remainder of this question, we will explore the density estimate for this data set. We will do so by generating a sample from the distribution with density function $f = \hat{f}$, the estimated density function produced in the previous plot. We will use rejection sampling to obtain this sample.

[Rejection sampling](#) transforms a sample from a distribution with density function $g(y)$ into a sample from a distribution with density function $f(y)$. Rejection sampling is useful when it is easy to sample from $g(y)$ but difficult to sample from $f(y)$. Rejection sampling uses a constant $M > 0$ such that $f(y) \leq Mg(y)$ for all y to repeatedly implement the following steps:

- Generate $y^* \sim g(y)$ and u^* from the uniform $\mathcal{U}(0, 1)$ distribution
- If $u^* \leq \frac{f(y^*)}{Mg(y^*)}$, keep y^* . Otherwise, reject y^* .

It can be shown that the sample of the retained y^* values is indeed a sample from $f(y)$. We illustrate this fact in the below figure with a simple example, where $g(y) = 1$ for $0 < y < 1$ and $f(y) = 2y$ for $0 < y < 1$. Here, we used $M = 2$. The left plot shows the (y^*, u^*) combinations that were retained in green and those that were rejected in red. When considering the red and green points together, we can see that y^* values were originally generated from a uniform $\mathcal{U}(0, 1)$ distribution that corresponds to $g(y) = 1$ for $0 < y < 1$. The

right plot shows the histogram of the y^* values that were retained. This histogram aligns nicely with the density function $f(y) = 2y$ for $0 < y < 1$ given by the blue line in the left plot.



You will implement rejection sampling in stages over the next parts of this question.

- (d) [2 points] First, you will generate 50000 y^* values from the uniform $\mathcal{U}(-1, 18)$ distribution. This can be done using built-in functions in **base R**. Then, compute $\hat{f}(y^*)$ for each of these values using the `kde_gaussian()` function you created in Question 3 with the bandwidth `h` you computed in part (c) of this question.
- (e) [3 points] Next, you will generate 50000 u^* values. You will then compute the ratios $f(y^*) \div M g(y^*)$ from the second bullet point earlier to decide whether to reject or retain the y^* value. You should use $M = 7.25$ and discern f and g based on information provided earlier in this question.
- (f) [2 points] You will now plot a histogram of your retained y^* values on the relative scale, where the bins for the histogram are unshaded and created using `breaks = "fd"`. The limits for the x -axis of this plot should coincide with `xrange` used in part (c) of this question. Remember to include informative titles and axis labels.
- (g) [3 points] Lastly, you will use the sample of your retained y^* values to estimate $Pr(Y < 1)$, where Y is the random variable with density function \hat{f} from part (c). Based on the formula for the kernel density estimator from Question 2, why is this probability greater than 0 when the Gaussian kernel is used? Briefly, what do the results from this question suggest about using kernel density estimation to summarize data from distributions with bounded support.

B | STAT 341 Assignment Walkthrough Videos

Provided below are Google Drive links to the four assignment walkthrough videos that I created for STAT 341. These videos were initially posted on the University of Waterloo's learning management system. These videos clarified the expectations and rationale for each assignment question. Students were required to submit their answers as PDFs that were generated by R Markdown. The first video in particular clarifies the formatting requirements. While these videos were explicitly classified as optional viewing, approximately 75% of the students in my section and the other STAT 341 section viewed at least one of the videos.

I set out to create these videos for several reasons. First, I thought filming these videos might preempt some simpler questions from being asked repeatedly on our course discussion board, which could help save the instructional team time throughout the term. Second, several of my students shared that they mainly worked on assignments in the evenings or on weekends. I did not expect the instructional team to be answering questions at those times; these videos provided an asynchronous source of support when the instructional team was not available. Third, I wanted to gain practice recording instructional videos given the recent increased prevalence of online and hybrid courses. If necessary, I feel more prepared to teach such courses after recording these videos.

- [Google Drive link for Assignment 1 Video](#)
- [Google Drive link for Assignment 2 Video](#)
- [Google Drive link for Assignment 3 Video](#)
- [Google Drive link for Assignment 4 Video](#)

C | STAT 938 Mock Consulting Session Plans

STAT 938 Mock Consulting Sessions – Spring 2023

Background: I am a member of the UW Graduate Student Association. We want to provide students with guidance to see if they are paying a fair price for rent. We hope to use a linear regression model to provide this guidance. I am a graduate student in urban planning. I took a statistics course in my undergraduate degree that covered basic regression, but it has been a while. Data from 190 (mock) apartment listings on rentals.ca from Waterloo region have been collected (simulated). I will say these were all listings available on the website when I collected data last week. The response variable is rent price (in \$CAD/month).

The model contains the following variables (bold factors are significant):

- Binary: **Utilities included**, **Laundry included**, Air conditioning included.
- Categorical: Cambridge vs. Waterloo and Kitchener vs. Waterloo
- Continuous: **Number of bedrooms**, **Age of building** (years), **Square feet**, Number of bathrooms

Questions: Each student will get a “warm-up” question (interpreting a regression coefficient) and a “follow-up” question.

Warm-up questions (i.e., we are particularly interested in understanding what this variable means):

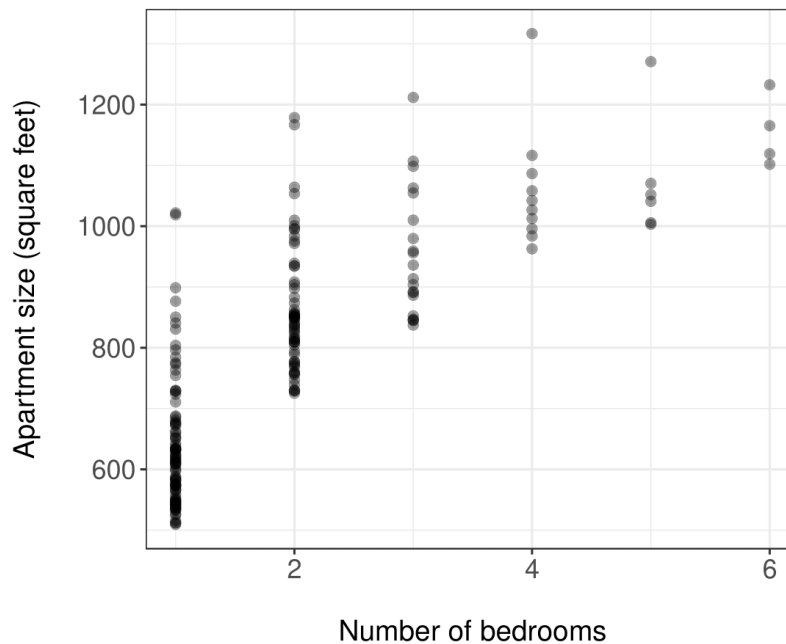
- *Question 1: Interpret utilities variable.* For two similar apartments, we expect the monthly rent to be about \$183 more expensive when utilities are included. Similar means that they both have the same number of bedrooms, city, etc. as in all other factors in the model do not change.
- *Question 2: Interpret the laundry variable.* For two similar apartments, we expect the monthly rent to be about \$91 more expensive when utilities are included.
- *Question 3: Interpret number of bedrooms variable.* For two similar apartments, we expect the monthly rent to increase by about \$119 for each bedroom that is added to the apartment.
- *Question 4: Interpret the building age coefficient.* For two similar apartments, we expect the monthly rent to decrease by about \$3 for each year increase in the age of the apartment building.

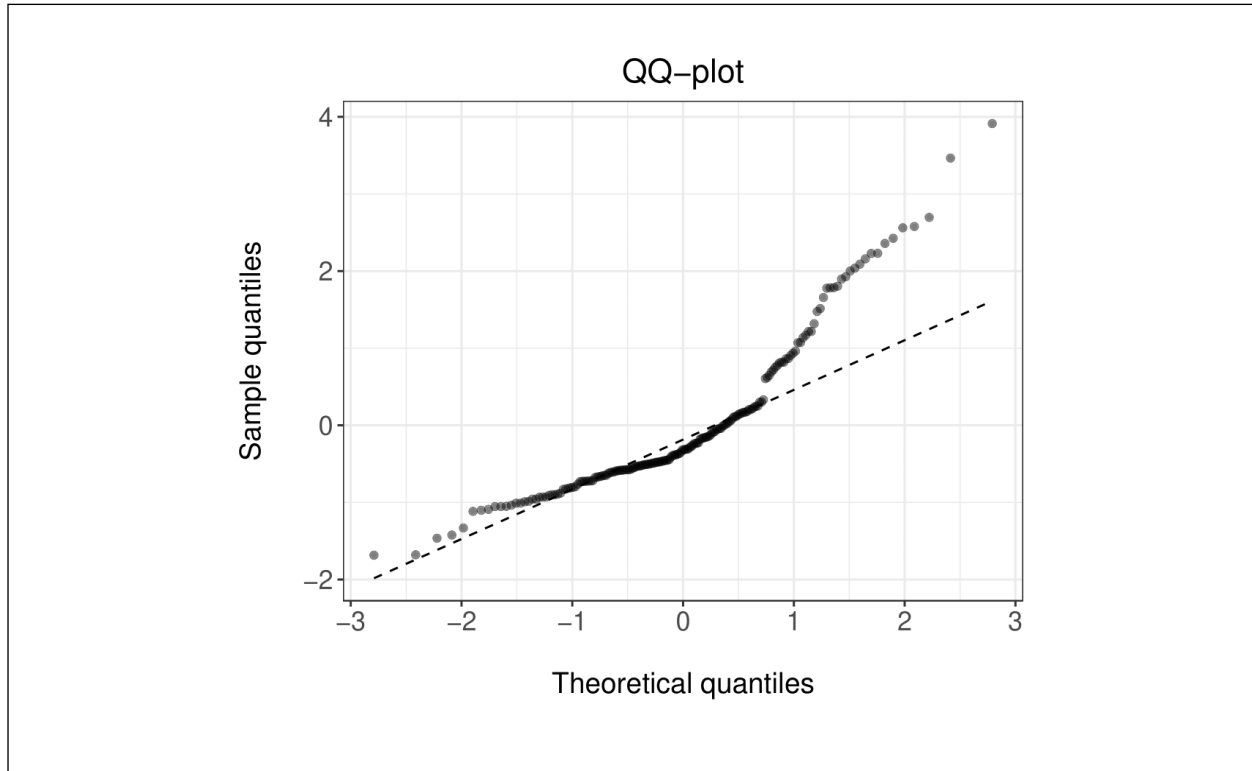
Follow-up questions:

- *Question 1:* I just put all the apartment information that I collected into the model because I didn’t want to overlook the impact that any of them might have on rent. Do you have advice for how I should choose which factors to put into the model?
 - Students should mention something about variable selection (i.e., forward, backward, stepwise selection rules, cross-validation, regularization methods, etc.).
- *Question 2:* I was investigating the relationship between some of the factors in the model and produced this plot. It seems like the apartment size is related to the number of bedrooms. Is it okay if the factors in my model are related to each other? How should I proceed?
 - Students should mention something about multicollinearity and discuss variance inflation factors. They may also discuss variable selection techniques.
- *Question 3:* I was looking through my statistics notes from a while ago, and they mentioned something about a QQ-plot. What does this plot tell me? Does it look okay for this model? Do you have any suggestions for how to fix this?
 - Students should mention that the residuals do not appear to be normal (they are very skewed). They should mention transforming the response variable as a solution (i.e., perhaps take the response for the model to be the logarithm or square root of rent).

Model Output and Plots

```
##
## Call:
## lm(formula = Rent ~ UtilitiesIncluded + factor(City) + LaundryIncluded +
##     AirConditioningIncluded + NumberOfBedrooms + BuildingAge +
##     SquareFeet + NumberOfBathrooms, data = apartment)
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      790.69    126.58      6.2   3e-09 ***
## UtilitiesIncluded    183.04     42.04      4.4   2e-05 ***
## factor(City)Cambridge  -53.47     53.69     -1.0   0.321
## factor(City)Kitchener  -39.99     58.62     -0.7   0.496
## LaundryIncluded      90.75     41.19      2.2   0.029 *
## AirConditioningIncluded   9.51     41.27      0.2   0.818
## NumberOfBedrooms     118.86     30.75      3.9   2e-04 ***
## BuildingAge         -3.16       0.96     -3.3   0.001 **
## SquareFeet           1.05       0.18      5.9   1e-08 ***
## NumberOfBathrooms      64.62     36.06      1.8   0.075 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 280 on 180 degrees of freedom
## Multiple R-squared:  0.64, Adjusted R-squared:  0.63
## F-statistic: 36 on 9 and 180 DF, p-value: <2e-16
```





The intended learning outcomes for this activity were to

- Apply effective questioning techniques to understand the client's problem prior to offering statistical guidance.
- Communicate a statistically valid solution to the client in a manner that is accessible to non-statisticians.

Each student had 10 minutes to conduct their meeting with the mock client (me). A rough plan for each 10-minute block is provided below.

| Time | Activity |
|-------------|---|
| 30 SECONDS | Mock introduction between client and consultant. This would normally take longer, but a brief and professional introduction will suffice here. |
| 2.5 MINUTES | Consultants will ask me to describe my problem and ask questions for confirmation. I will show them the regression model and ask them a warm-up question. |
| 2 MINUTES | Consultant provides an answer to the client's question with nontechnical jargon. |
| 2.5 MINUTES | Client will ask one of the follow-up questions now that consultants are more familiar with the model. |
| 2 MINUTES | Consultant comments on the follow-up question with nontechnical jargon. |
| 30 SECONDS | Briefly wrap up the meeting. Consultants may ask the client to reach out by email or book another appointment if they have questions. |

D | BIOS 612 Jupyter Notebooks

Provided below are Google Drive links to 12 Jupyter notebooks that I created for BIOS 612. These notebooks were originally posted on McGill University's learning management system as dynamic files that students could engage with during class. The links below prompt downloads of static HTML versions of the notebooks that were completed by the end of class.

- [Notebook for Bayesian treatment of missing data](#)
- [Notebook for frequentist treatment of missing cross-sectional data](#)
- [Notebook for frequentist treatment of missing longitudinal data](#)
- [Notebook for generalized mixed models with conditional likelihood](#)
- [Notebook for generalized mixed models with marginal likelihood](#)
- [Notebook for Bayesian generalized mixed models](#)
- [Notebook for generalized estimating equations](#)
- [Notebook for exploratory analysis of correlated binary data](#)
- [Notebook for frequentist modeling of correlated binary data](#)
- [Notebook for Bayesian modeling of correlated binary data](#)
- [Notebook for frequentist spline methods](#)
- [Notebook for Bayesian spline methods](#)

I provide some context for these notebooks using excerpts from the notebook for generalized mixed models with conditional likelihood. In the notebook, students are first introduced to a real data set that assesses the efficacy of a treatment for seizures.

Seizure Example - Conditional GLMMs

We now consider a log-linear model for the seizure data, one of the examples from Chapter 9 of the Wakefield textbook.

For each of 59 epilepsy patients, the no. of seizures were recorded during a baseline period of 8 weeks, after which patients were randomized to treatment with the drug progabide or to placebo. The no. of seizures was then recorded in 4 consecutive 2-week periods. Patient age was also available. Let

$$\begin{aligned} Y_{ij} &= \text{number of seizures on patient } i \text{ at occasion } j \\ t_{ij} &= \text{length of observation period on patient } i \text{ at occasion } j \\ x_{i1} &= 0/1 \text{ if patient } i \text{ was assigned placebo/progabide} \\ x_{ij2} &= 0/1 \text{ if } j = 0/1, 2, 3, 4 \end{aligned}$$

with $t_{ij} = 8$ if $j = 0$ and $t_{ij} = 2$ if $j \geq 1$, for all i .

```
suppressPackageStartupMessages(require(ggplot2)) # Load packages
suppressPackageStartupMessages(require(dplyr))

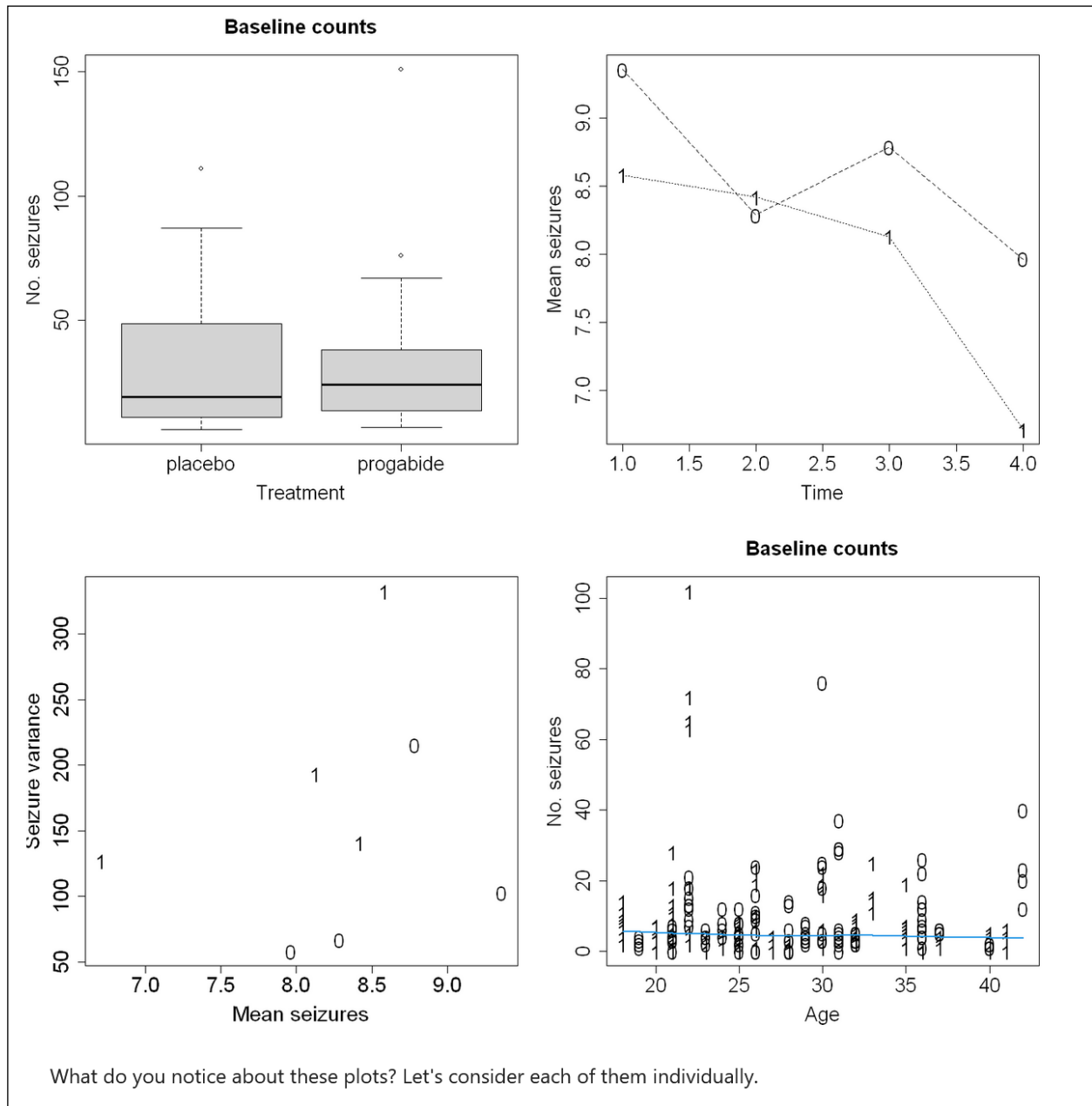
seiz.dat <- read.csv("TV1998.csv", header=T) # Load and format data
head(seiz.dat)
```

A data.frame: 6 × 9

| | y | trt | base | age | V4 | subject | period | lbase | lage |
|---|-------|---------|-------|-------|-------|---------|--------|------------|------------|
| | <int> | <chr> | <int> | <int> | <int> | <int> | <int> | <dbl> | <dbl> |
| 1 | 5 | placebo | 11 | 31 | 0 | 1 | 1 | -0.7563538 | 0.11420369 |
| 2 | 3 | placebo | 11 | 31 | 0 | 1 | 2 | -0.7563538 | 0.11420369 |
| 3 | 3 | placebo | 11 | 31 | 0 | 1 | 3 | -0.7563538 | 0.11420369 |

Students were then provided with code to produce the plots on the next page. In class, we discussed the students' takeaways from these plots. For instance, we can observe in the top right

plot that the progabide treatment (indicated by 1s) might be more effective than the placebo (indicated by 0s) at reducing seizure counts. More importantly, the bottom left plot reveals that the variance of the seizure counts is much greater than the average seizure count.



Given the students' existing knowledge of the mean-variance relationship for count data modeled using a Poisson distribution, they should be able to identify that there is substantial overdispersion in the data. This exercise led into a discussion about how we can account for this overdispersion for correlated count data in a conditional likelihood framework. Ultimately, the notebooks in this course empowered students to develop their own intuition for drawing inference from complex data.

E | R Shiny Application for CUT Program

This [app](#) was used to facilitate the following interactive activity in my CUT project.

Interactive Activity: Active Learning Strategies with p -Values

Background: We use an interactive application to explore some features of the p -value using simulation: https://lmhagar.shinyapps.io/CUT_App/ We use a mock experiment to assess the impact of active learning in the classroom. We are facilitating two versions of the same course, one with traditional lectures and the other with flipped lectures. We compare the effectiveness of the two lecture formats using the students' final exam marks. To do this, we test the hypothesis that the average exam marks in the two versions of the course are the same.

Instructions: The simulation settings can be controlled using the *Simulation Inputs* tab. We simulate an exam mark for each student in the experiment, such that the average mark in the traditional lecture is 70. The simulation has two inputs:

- *Average mark in flipped lecture*: marks in the flipped lecture are simulated such that the average mark is determined by this slider. It can take whole numbers between 64 and 76.
- *Number of students per class*: this input controls how many marks we simulate for *each* class. It can be increased from 5 to 125 in increments of 20.

The active learning experiment with these settings is repeated 1000 times each time the *Generate* button is pressed. For each repetition, we obtain a p -value. We use the histogram of these 1000 p -values to explore features of the p -value. Each p -value tells us how compatible the simulated data are with the hypothesis of the two average marks being the same. The red line on the histogram shows a 5% cutoff. This cutoff is typically used in practice (i.e., the hypothesis of the two average marks being equal is rejected if the p -value is less than 5%). The plot title conveys the proportion of the 1000 simulations in which the p -value was less than 5%. The simulation can be rerun adjusting the slider inputs (if necessary) and pressing the *Generate* button.

Activity 1: We start with the default settings (average mark in flipped lecture = 70, number of students per class = 5).

- Please complete this part **before** pressing *Generate*. For this setting, the hypothesis is true. The simulation will generate a histogram of p -values. What do you think the histogram will look like?
- Repeat the simulation a few times with these settings. What patterns do you see in the histogram? Is this what you expected to see? Approximately what proportion of the p -values are less than the 5% cutoff?

Please pause here.

- Please move the *Number of students per class* slider to a larger sample size. **Before** you press *Generate*, how do you think the histogram of p -values will change from those for the initial setting?
- Repeat the simulation a few times with these settings. What changes to the histogram do you notice? Is this what you expected to see? Approximately what proportion of the p -values are less than the 5% cutoff?

Please pause here.

- e) What proportion of p -values do you think would be less than a 10% cutoff for this scenario?
- f) For this simulation setting, would rejecting the hypothesis be the correct decision? Why might 5% be a popular cutoff value then?

Please pause here.

Activity 2: Please return the *Number of students per class* slider to 5 students per class. Think of a difference d between the two average marks that would be of practical importance **to you**. For this simulation, d should be one of 1, 2, 3, 4, 5, or 6. For instance, if the difference between the two averages is less than d , we can consider the averages for the two classes to be practically the same. Otherwise, there is a notable difference in the average marks. Move your average mark in flipped lecture to $70 + d$.

- a) Please complete this part **before** pressing *Generate*. For this setting, what do you think the histogram of p -values will look like?
- b) Repeat the simulation a few times with these settings. What patterns do you see in the histogram? Is this what you expected to see? Approximately what proportion of the p -values are less than the 5% cutoff?

Please pause here.

- c) Please move the *Number of students per class* slider to 25. **Before** you press *Generate*, how do you think the histogram of p -values will change from those from part a)?
- d) Repeat the simulation a few times with these settings. You could also try this with larger sample sizes. What changes to the histogram do you notice? Is this what you expected to see? Approximately what proportion of the p -values are less than the 5% cutoff?
- e) For this simulation setting, what is the correct decision with regard to rejecting or not rejecting the hypothesis? Are you likely to make the correct decision if you have a small sample size of students?

Activity complete!