

Methods of Data Analysis II

University of Toronto
Department of Statistical Sciences
STA303H1S Winter 2024

	LEC 0101:	LEC0201:	LEC5101:
Monday Lectures	9am-11am ET at AH100	11am-1pm ET at AH100	5pm-7pm ET at AH100
Wednesday Lectures	9am-10am ET at AH100	12pm-1pm ET at AH100	2pm-3pm ET at AH100
Instructor OH:	W 10:40am-11:30pm ET at	M 2:15pm-3pm ET at MP	W 3:30pm-4:30pm ET
OH Locations:	Health Sciences Room: 108	McLennan Physical Laboratory: 118	Multiple Locations

Office Hours will start on the third week. There will be additional office hours held by the instructors before the final project and the final exam.

COURSE OVERVIEW

Course Description: Briefly, the contents will focus on the categorical data analysis, exponential family, generalized linear models (GLMs), quasi-likelihood, generalized estimating equations (GEEs), linear mixed models, generalized linear mixed models (GLMMs) and generalized additive models (GAMs). Statistical analysis will be conducted using R. The students will be expected to simulate datasets using R, as well as interpret R code and output on tests and projects.

Content, emphasis, etc. of the course is defined by means of the lecture material — not only the posted lecture notes. It is important to attend all lectures, as there is normally no simple way to make up for missed lectures (perhaps obtain another student's notes). There will also be many lecture examples using statistical software R.

Learning Outcomes: By the end of this course, all students should have a solid understanding of both the mathematical theory of GLMs, as well as their application in the form of data analysis. The course will also focus on the applications of LMMs, GLMMs and GAMs in real life scenarios. Students should be prepared to show their understanding of the above through:

- application of methods through problem-solving questions;
- description and explanation of concepts relating to the mathematical theory;
- derivation and proof of topics based on GLM concepts and theory;
- practical application of methods on real data using statistical software R, with appropriate justification of use of these methods;
- interpretation of data analysis results in clear and non-technical language

Pre-requisites: The prerequisites are:

Prerequisite: STA302H1 or STAC67H3 or STA302H5

Exclusion: STAC51H3

Pre-requisites are **strictly enforced by the department, not the instructor**. If you do not have the equivalent pre-requisites, you will be un-enrolled from the course.

COURSE MATERIALS

Course Content: All lecture slides, recordings and materials will be posted on the Quercus course page for each lecture section. Furthermore, any important announcements will also be posted in Quercus. Please make sure to check it regularly so you don't miss anything.

Textbook: There are no required textbooks. All assessments will be conducted based on lecture materials. However, there are some useful books which will be used for references.

1. Alan Agresti. “*Categorical Data Analysis (3rd edition)*”. Wiley, 2011.
2. David Clayton and Michael Hills. “*Statistical Models in Epidemiology*”. OUP Oxford, 2013.
3. Peter McCullagh and John Ashworth Nelder. “*Generalized Linear Models (2nd Edition)*”, Chapman & Hall/CRC, 1989.
4. Peter Diggle, Patrick Heagerty, Kung-Yee Liang, Scott Zeger. “*Analysis of Longitudinal Data*”. Oxford University Press; 2002
5. Julian Faraway. “*Extending the Linear Model with R: Generalized Linear, Mixed Effects and Non-parametric Regression Models*”. CRC press, 2016.
6. Simon N. Wood. “*Generalized Additive Models: An Introduction with R*”. CRC press, 2017.

Statistical Software: We will be using the R Statistical Software for performing statistical analyses in this course. R is a free software that can either be downloaded onto your personal computer or used in a cloud environment. We encourage all students to use RStudio through the [JupyterHub](#) for University of Toronto. This will allow you to login with your official UofT credentials and use RStudio without the need for a local installation and can be run on any device that has access to an internet connection. More information about using RStudio in JupyterHub will be provided early in the term. R code shown in class will be available on the course page and, along with any additional resources, should be sufficient to complete any assessment involving data analysis.

COURSE COMPONENTS

Lectures: Lectures will be conducted in person in [AH100](#). Slides will be available after the class. Class time each week will comprise of a combination of lecturing, and code-along sessions. Where possible, you are encouraged to bring a laptop or tablet to follow along with the code.

Office Hours: Instructor and TAs will hold office hours in a combination of online and in-person formats. The office hour schedule and mode of delivery will be posted on Quercus once finalized. It is recommended that you visit office hours whenever you have a question about the material. It is always important to have material clarified as quickly as possible. Don't wait until the last minute to ask your questions!

Piazza: We will be using the Piazza as an online discussion forum, which can be accessed through the Quercus course page. **All questions about course material should be posted here** or asked during TA/instructor office hours. The instructor and TAs will monitor the board and will help answer questions but students are encouraged to answer posts and help their fellow classmates.

COMMUNICATION

How your instructor will communicate with you: All communication will be made through Quercus announcements or during lectures. Please ensure that you check Quercus regularly so you don't miss anything important.

Where to send content questions: We will be using the Piazza to collect student questions regarding course content, assignments, etc. All questions should be posted here.

When to email the instructor: The instructor will only respond to emails of a private or sensitive nature. If you email the instructor with content related questions, you will be asked to repost your question on the content board so the answer may benefit all students. Should you need to email the instructor about a sensitive or personal nature, please use your official mail.utoronto.ca email, include your full name and student number in the text. Send all course related emails to sta303@utoronto.ca. Please allow up to 48-96 hours for a reply. Emails will not be monitored on evenings and weekends.

A note on email and discussion board etiquette: Please make sure that you communicate politely and respectfully with all members of the teaching team and your fellow classmates. Written communications can sometimes take a tone other than what was intended (e.g. can come off as dismissive, rude or insulting), so make sure you re-read or read out loud your email/post before sending it to make sure it has the tone you intended. For more tips on respectful communication, see [professional communication tips](#). Piazza is a teaching and learning tool and therefore should only be used as such. Any posts that detract from the learning goal of the board will be removed to keep the board a safe space.

GRADING SCHEME

All the students will be evaluated in the following way:

Assessment	Date	Weight
Term Test	February 12	25%
Assignment	February 26	10%
Final Project Proposal	March 11	10%
Final Project Report	April 3	25%
Final Exam	April 10-30	30%

Please note that the last day to drop the course without penalty is March 11, 2024.

EVALUATION BREAKDOWN

Term Test: The term test will be conducted in person during the scheduled Monday class time (see top of page 1). The test will be approximately 2 hours long. More details will be communicated closer to the test date. The test will cover material from Weeks 1-5.

Assignment: You will be given one assignment in the term. The purpose of this assignment is to develop your understanding of the statistical properties of the estimators obtained from a generalized linear models. This will be useful for developing data analysis skills as well as to develop practical understanding of the methods taught in the class. The assignment will have a heavy focus on the use of statistical software (R specifically), and will involve applying the methods learned during lecture to a data set. The format of the assignments will be as follows:

1. use the methods taught in lecture to perform a small data analysis.
2. simulate unique datasets and writing your own functions instead of built in R functions.

3. solve some mathematical problems and explain the procedure with simulated datasets

Final Project: The final project will be due on the last day of the lectures and will consist of a data analysis on **a novel dataset** of your choice. Students will be required to demonstrate their understanding of the methods taught in lecture by developing a reasonable GLM/GLMM/LMM model that addresses a valid research question using the techniques taught in class. The students will be responsible for choosing the correct methods to apply and providing appropriate justifications defending their choices. The final project is a scaffolded assessment involving 2 parts:

- Part 1- Research question and dataset selection: Students must find a dataset available online and define a research question that can be answered with this dataset using GLM. Students will need to explain why their research question is important and how GLM/GLMM/GAMM may be used to answer it. A short exploratory data analysis of the chosen dataset will also be required. Students will then be required to prepare a **five minute long video presentation**, which have to be uploaded on Quercus. More details will be provided during the lectures. This part will be due on March 11th.
- Part 2 - Final Project Report: Students will put together a scientific report that outlines the relevance of their proposed research question, the process of their analysis, the results of the performed data analysis, and a discussion of the meaning of the results as well as limitations of the analysis with respect to the statistical tools used/decisions made or the data used. This part will be due on April 3rd (the last day of the lectures).

The final project will be done individually, and must be typed and submitted by the deadline. More detailed instructions will be provided at a later date.

Final Exam: The details about the final exam will be provided during the last week lectures. For the final exam we will be following standard University of Toronto Schedule. the final exam will be three hours in duration and will be scheduled by the Faculty of Arts and Science during the final assessment period.

LATE ASSESSMENT AND EXTENSION REQUEST POLICY

The assessment deadlines may change from the ones stated in the syllabus depending on how the lecture progresses. However, once the deadline(s) has been announced, the students need to submit the assignments by the deadline. Students will be able to still submit the assignments up to 5 days after the deadline, however, each additional day will be accounted for 20% penalty.

Extreme Situations/Prolonged Illness Extensions: Should a student be experiencing a prolonged illness or other situation that prevents them from turning in their work by the deadline, they should **immediately contact their instructor and College Registrar** to inform them of their situation. They should also submit an **Absence Declaration form on ACORN** that lists every day during which they were incapacitated and unable to work. Accommodations or further extensions will not be considered without a completed declaration, and will only be considered for extreme circumstances.

Accessibility-Related Extension Requests: Students registered with Accessibility Services should notify the instructor as soon as possible if additional time is needed on assessments that are eligible for extensions. Please **notify the instructor by email of your situation and cc your accessibility advisor** in the process. The instructor will work with the accessibility advisor to determine an appropriate extension for your situation.

MISSED ASSESSMENT POLICY

If you experience a prolonged absence due to illness or emergency that prevents you from completing any number of assessments, please contact your College Registrar as soon as possible so that any necessary arrangements can be made.

Missed Assignment or Final Project: Missing assessments will receive a 0.

Missed Term Test: If a student is experiencing a serious personal illness or emergency on the date of the test, the student **must declare their absence on ACORN and notify the teaching team via email no later than one week after the date of the test.** A make-up test will then be scheduled at a date and time determined by the instructor. **The format of the make-up is at the discretion of the instructor and may not resemble the format of the original (e.g. an oral exam).**

REGRADE REQUESTS

Regrade requests will be accepted for all assessments. Regrade requests must provide a justification for where there exists a grading error and/or how the work meets the grading rubric. These justifications must further be backed up with concrete references to the course material. All regrade requests will be accepted through a form available on the Quercus course page and will be accepted no later than one week after the grade for that assessment is released. **No regrade requests will be accepted by email or after the 1 week deadline.** The instructor further reserves the right to re-evaluate the assessment in its entirety (i.e. grades can go up, down, or remain unchanged). Please allow a few weeks for regrade requests to be processed by the instructor.

INTELLECTUAL PROPERTY

Course materials provided on Quercus, such as lecture slides, assignments, tests and solutions are the intellectual property of your instructor and are for the use of students currently enrolled in this course only. **Providing course materials to any person or company outside of the course is unauthorized use.** This includes providing materials to predatory tutoring companies.

ACADEMIC INTEGRITY

The University treats cases of plagiarism and cheating very seriously. It is the students' responsibility for knowing the content of the University of Toronto's [Code of Behaviour on Academic Matters](#). All suspected cases of academic dishonesty will be investigated following procedures outlined in the above document. If you have questions or concerns about what constitutes appropriate academic behaviour or appropriate research and citation methods, you are expected to seek out additional information on academic integrity from your instructor or from other institutional resources (see <http://academicintegrity.utoronto.ca/>). Here are a few guidelines regarding academic integrity:

- Sharing or discussing questions or answers with other students during tests is an academic offence.
- Students must complete all assessments individually. Working together is not allowed.
- Paying anyone else to complete your assessments for you is academic misconduct.
- Sharing your answers/work/code with others is academic misconduct.
- All work that you submit must be your own! You must not copy mathematical derivations, computer output and input, or written answers from anyone or anywhere else. Unacknowledged copying or unauthorized collaboration will lead to severe disciplinary action, beginning with an automatic grade

of zero for all involved and escalating from there. Please read the UofT Policy on Cheating and Plagiarism, and don't plagiarize.

RULES REGARDING THE USE OF GENERATIVE AI IN ASSESSMENTS

Generative Artificial Intelligence (AI), and specifically foundational models that can create writing, Computer code and/or images using minimal human prompting are proliferating and becoming ubiquitous. This includes not only GPT-4 (and its siblings ChatGPT and Bing), but many writing assistants that are built on this or similar AI technologies. There are now hundreds of these systems that are readily available. In this course, the use of such AI tools is limited. The students are most welcome to use the tools for learning purposes and to gather information about the topics that are taught during the course. However, these tools should not be used to complete any of the assessments in the course. For details please read the following statements:

- Students may use artificial intelligence tools for creating an outline for an assignment, but the final submitted assignment must be original work produced by the individual student alone.
- Students may not use artificial intelligence tools for taking tests, writing project report, creating R codes, or completing major course assignments. However, these tools may be useful when gathering information from across sources and assimilating it for understanding.
- Representing as one's own an idea, or expression of an idea, that was AI-generated may be considered an academic offense in this course.
- Students may not copy or paraphrase from any generative artificial intelligence applications, including ChatGPT and other AI writing and coding assistants, for the purpose of completing assignments in this course.
- This course policy is designed to promote your learning and intellectual development and to help you reach course learning outcomes.

ACCESSIBILITY NEEDS

Students with diverse learning styles and needs are welcome in this course. If you have an acute or ongoing disability issue or accommodation need, you should register with Accessibility Services (AS) at the beginning of the academic year by visiting <https://studentlife.utoronto.ca/department/accessibility-services/>. Without registration, you will not be able to verify your situation with your instructors, and instructors will not be advised about your accommodation needs. AS will assess your situation, develop an accommodation plan with you, and support you in requesting accommodation for your course work. Remember that the process of accommodation is private: AS will not share details of your needs or condition with any instructor, and your instructors will not reveal that you are registered with AS.

CLASS SCHEDULE - TENTATIVE

This is the tentative outline for Winter 2024. Topics may be reduced or additional topics may be added by course instructor's discretion.

Week	Content
1 (Jan 8-10)	Review of Maximum Likelihood Estimation (MLE), Score, Wald, Likelihood Ratio based confidence intervals and test of hypothesis.
2 (Jan 15-17)	Study designs, contingency tables, risk difference and risk ratio. Odds ratio, rate ratio. Test of independence.
3 (Jan 22-24)	Small sample inference. Fisher's exact test. Delta method, confounding and interaction. Review on linear regression. Introduction to GLMs. Logistic regression.
4 (Jan 29-31)	Exponential family and iteratively reweighted least squares (IRLS). Poisson regression.
5 (Feb 5-7)	Model diagnostics. Variable selection: Stepwise methods, LASSO. Classification and Discrimination using Logistic Regression. Cross validation, ROC curve and Bootstrap method.
(Feb 12)	Term Test
6 (Feb 14)	Analyzing real life datasets with R
(Feb 19-23)	Reading Week
(Feb 26)	Assignment Due
7 (Feb 26-28)	Overdispersion and negative binomial regression. Offset terms in GLM. Multinomial GLM, proportional odds GLMs and probit GLM.
8 (Mar 4-6)	Matched Case-Control studies and conditional logistic regression
(Mar 11)	Final Project Proposal Due and Deadline to drop course without penalty
9 (Mar 11-13)	Survival analysis. Kaplan-Meier estimates. The Cox proportional hazards regression.
10 (Mar 18-20)	Linear mixed effects models (LMM). Generalized linear mixed effects model (GLMM) Quasi Likelihood and generalized estimating equation (GEE)
11 (Mar 25-27)	Scatterplot smoothing, Cubic Splines, Generalized Additive Models (GAM), Generalized Additive Mixed effects Model (GAMM)
12 (Apr 1-3)	Regression & Classification trees. Bagging and Random Forests.
(Apr 3)	Final Project Report Due
April 10-30	Final assessment period