

STA302: Methods of Data Analysis I Summer 2020

Nnenna Asidiana

E-mail: nnenna.asidiana@mail.utoronto.ca

Office Hours: TBA

Class Hours: Mon 6:00pm-9:00pm

Wed 6:00pm-9:00pm

- *This is an online course. Please note that since lectures and/or evaluations will be taking place during the above lecture times, you must be available during those times. No accommodations will be made for assessments missed during these times.*
- *As this is an online course and all assessments must be submitted through Crowdmark or Quercus, it is the student's responsibility to ensure they have a reliable internet connection.*

Course Description

This course covers theory and applications of regression analysis. We will develop the theory of regression models and study how to analyse data when such models are appropriate. Topics to be covered include: simple and multiple regression models using least squares, analysis of variance, inference for regression parameters when the errors are normally distributed, confidence and prediction intervals, geometry of least squares, multicollinearity, regression models for quantitative and qualitative predictors, model selection and validation, and diagnostics.

Prerequisites/Corequisites

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. Students should have a second year statistics course, such as STA238, STA248, STA255, or STA261, a computer science such as CSC108, CSC120, CSC121, or CSC148 and a mathematics course such as MAT221(70%), MAT223, or MAT240 or equivalent preparation as determined by the department.

Email Policy

Important announcements, lecture notes, additional material, and other course info will be posted on Quercus. Check it regularly. You are responsible for keeping up with announcements from instructors on Quercus and via e-mail.

To ensure your email gets to me, please ensure the following:

- Use your academic email, for example nnenna.asidanyal@mail.utoronto.ca
- Use the following format in the subject of your email: CourseName/LASNTNAME For example: STA302S/ASIDIANYA
- Be very clear and concise.

Required Materials

- This course requires the following textbooks:
 1. Kutner, M. H., Nachtstein, C. J., and Neter, J. **Applied Linear Regression Models.** McGraw-Hill, 5th edition.
 2. Sheather, S.J. **A Modern Approach to Regression with R.** (Springer).
- We will be using RStudio for performing statistical analyses. R is a free software that can either be downloaded onto your personal computer or used in the cloud. If you choose to work with R on your personal computer, then installation will be a two step process:
 1. The base R framework is available for download at [R framework](#)
 2. RStudio is a good integrated development environment to R (makes it simpler to work in R) and can also be downloaded for free at [RStudio](#).

If you don't want to download the program or run into problems with installation, you may want to consider [rstudio.cloud](#) **RStudio Cloud** only requires you to login with your Utoronto email and connect to our course project via the link provided.

Course Components

Class Lectures: Two three hours a week class times will be used to cover important course materials. It is important that you attend these classes in order to keep up with the topics, and gain a deeper understanding of the applications of statistics in social sciences.

Office Hours: We will hold office through Bb Collaborate in the Quercus course page. The office hour schedule will be posted on Quercus. It is recommended that you visit office hours whenever you have a question about the material. This is very important since this is an online accelerated class, so onus is up to the student to have material clarified as quickly as possible. Don't wait until the last minute to ask your questions.

Course Assessment

Assignments: Over the course of the semester students will be given a mini assignment to work on that is related to the course material up until that point. This may include both an R component and a written component. More details, such as the content and deadline, will be communicated later. **No late report will be accepted.**

Your final grade for the course will automatically be determined by the higher of the two following grading schemes:

Undergraduate students will be evaluated in the following way:

Item	Date Due	% of Grade
Weekly quizzes (x5)	Open at the end of each Wednesday lecture	40%
Term Test 1	Due July 22 by 12:00PM EST	15%
Mini Project # 1	Due July 29 by 11:59PM EST	20%
Final Project	Due Aug. 12 by 11:59PM EST	25%

Graduate students will be evaluated in the following way:

Item	Date Due	% of Grade
Weekly quizzes (x5)	Open at the end of each Wednesday lecture	35%
Term Test 1	Due July 22 by 12:00PM EST	15%
Mini Project # 1	Due July 29 by 11:59PM EST	25%
Final Project	Due Aug. 12 by 11:59PM EST	25%

There will be one hour at the end of each Wednesday lecture period to submit your online quiz via Quercus. It will be based on the material from the previous week. No late submissions will be accepted.

Weekly Quizzes

There will be 5 “weekly” online quizzes, that will be open during the last hour of each Wednesday lecture. Quizzes will begin on July 15 and continue until the last lecture period.

- We will take the best 3 quiz marks and drop the worst 2 in the calculation of your overall quiz mark
- The quizzes will be multiple choice and cover material from the previous set of lectures. You may wish to have a calculator available at this time to aid in any calculations.
- Quizzes can be found under Quercus Quizzes in the navigation bar. Quizzes must be done individually.

- Missed quiz: Because only the best 5 quiz marks will be counted, there will not be any accommodations for missed quizzes. These will receive a mark of 0, but will be dropped as part of the two worst quiz marks. Therefore, you may miss two quizzes without penalty.
- **There are no make-up quizzes.** Quizzes, beyond the 2 that will be dropped, will be given zero.

Term Test

There will be one term test that will be released July 22nd 8:00AM EST, and will be due 12:00PM EST July 22nd. The term test will cover all the material covered in class and exercises until the last week before the test. Information about the break down of the test will be released closer to the exam date.

Mini Project

You will be given one mini project in the term. The purpose of this mini project is to develop your data analysis skills which will be useful for the final project at the end of the term. The mini projects will have a heavy focus on the use of statistical software (R ally), and will involve applying the methods learned during lecture to a dataset.

Final Project

The final project will be due on Aug. 12, 2020 by 11:59PM EST and will consist of a data analysis on a unique dataset. Students will be required to demonstrate their understanding of the methods taught in lecture by developing a reasonable regression model using the techniques taught in class. The students will be responsible for choosing the correct methods to apply and providing appropriate justifications where necessary. **This is a formal report and therefore it must contain the following sections:**

- Introduction section: provides details regarding the question you wish to address, why the model is being developed, how you intend to go about developing the model, and finally how the model meets the purpose mentioned earlier.
- Exploratory data analysis section: a detailed description of the variables in the data set with appropriate tables or figures that highlight certain characteristics of your variables that you deem important to mention.
- Model development section: a detailed discussion of the process used to come to the final model. Justifications may be both statistical and empirical in nature. You should also have as well as in-depth diagnostics to illustrate the 'goodness' of the model.
- Conclusion section: restate why the model is useful in the context of the data, provide an interpretation of the final model in non-technical language (i.e explain how the variables work, discuss predictions), and discuss any limitations/problems remaining with the model and how they might impact its use in the real world.

The final project will be done in groups of 3 or 4, and must be typed and submitted by the stated deadline. A word count limit will be given, as well as other more detailed instructions at a later date. In order to pass the course, *everyone* must submit the final project. At the end of the submission include a detailed paragraph of each group members contribution to the project. **This will be used to determine individual final grades.**

Missed Assessment Policy

Students are responsible for completing all of the assessments detailed in the previous section. If a student is sick and needs to request an extension or accommodation on a mini project, they must send an email to their instructor. In order for the request to be considered, the email:

- must be received **at least** 24 hours before the mini project or term test is due
- the subject line must be written in the format shown in *Email Policy*
- must include your full name and student number in the body of the email
- must specify for which project the extension/accommodation is being requested
- must include the following sentences:
 - “I affirm that I am experiencing an illness or personal emergency and I understand that to falsely claim so is an offence under the Code of Behaviour on Academic Matters.”
 - “I understand that the weight of this assessment will be distributed across the weekly quizzes (10%/15%) and the final project (5%).”

Remark Policy

Any requests to have marked work re-evaluated must be made in writing **within one week of the date the work was returned to the class**. The request must contain a justification for consideration. You are responsible to check that your scores are entered correctly on Quercus. Any requests for a mark that was not entered correctly in Quercus must be made in writing within one week of the date the mark was entered in Quercus.

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.

Accessibility Statement

Students with diverse learning styles and needs are welcome in this course. The University of Toronto offers academic accommodations for students with disabilities. If you require accommodations, or have any accessibility concerns about the course, the classroom, or course materials, please contact Accessibility Services as soon as possible: accessibility.services@utoronto.ca or <http://accessibility.utoronto.ca>

Academic Integrity Statement

Academic integrity is essential to the pursuit of learning and scholarship in a university, and to ensuring that a degree from the University of Toronto is a strong signal of each student's individual academic achievement. As a result, the University treats cases of cheating and plagiarism very seriously. The University of Toronto's Code of Behaviour on Academic Matters (<http://www.governingcouncil.utoronto.ca/policies/behaveac.htm>) outlines the behaviours that constitute academic dishonesty and the processes for addressing academic offences. Potential offences include, but are not limited to:

IN PAPERS AND ASSIGNMENTS: Using someone else's ideas or words without appropriate acknowledgement. Submitting your own work in more than one course without the permission of the instructor. Making up sources or facts. Obtaining or providing unauthorized assistance on any assignment.

ON TESTS AND EXAMS: Using or possessing unauthorized aids. Looking at someone else's answers during an exam or test. Misrepresenting your identity.

IN ACADEMIC WORK: Falsifying institutional documents or grades. Falsifying or altering any documentation required by the University, including (but not limited to) doctor's notes. All suspected cases of academic dishonesty will be investigated following procedures outlined in the Code of Behaviour on Academic Matters. 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/>).

Schedule and weekly learning goals

Table 1: Tentative Weekly Schedule

Week 1	<p>Relation between variables (p.2-3)</p> <p>Regression models and their uses (p.5-9)</p> <p>Data for regression analysis (p.12-13)</p> <p>Overview of steps in regression analysis (p.13-15)</p>	Introduction
Week 1 cont.	<p>Model and assumptions (p.9-12)</p> <p>Least squares estimation (p.15-22)</p> <p>Predicted values, residuals (p.22-23)</p>	Simple linear regression I
Week cont.	<p>Analysis of variable in the dependent variable (p.63-68)</p> <p>Normal errors regression model, maximum likelihood estimation (p.2632)</p> <p>Precision of estimates (p.40-63)</p>	Simple linear regression II
Week 2	<p>Inference (p.40-63)</p> <p>Gauss-Markov theorem (p.18-19)</p> <p>Regression through the origin (p.161-165)</p> <p>Matrices and vectors (p.176-197)</p>	Simple linear regression III

Table 2: Tentative Weekly Schedule

	Multiple linear regression and matrix approach (p.197-199,214-223)	
Week 2 cont.	Estimation (p.199-201,223-224)	Multiple Linear Regression I
	Fitted values and residuals (p.202-204,224-225)	
	Properties of linear functions of random vectors	
Week 3	Properties of regression estimates (p.227-232)	Multiple Linear Regression II
	Normal errors, maximum likelihood estimation	
	Analysis of variance and quadratic forms (p.204-206,225-227)	
Week 3 cont.	Hypothesis testing and confidence intervals for single regression parameter, mean response, and prediction of a new observation (p.227-232)	Multiple linear regression III
	Hypothesis testing for several regression parameters: ANOVA and distribution of quadratic forms (p.256-268)	
	Multicollinearity (p.278-289)	
Week 4	The geometry of least squares Violation of assumptions	Multiple linear regression IV

Table 3: Tentative Weekly Schedule

Week 4	cont.	An example with R	Interlude
Week 5		Polynomial regression model (p.294-305) Interactions (p.306-313) Qualitative predictors (p.313-321)	Regression models for quantitative and qualitative predictors
Week 5 cont.		Overview (p.343-350) Analyzing Associations Between Categorical Variables Criteria for model selection (p.353-361) Model validation (p.369-375)	Model selection and validation
Week 6		Model adequacy (p.384-390) Outlying Y observations (p.390-398) Outlying X observations (p.398-400) Influential cases (p.400-406)	Diagnostics