

QBIO 310: Statistical Thinking for Quantitative Biology Syllabus

General Information

Lecture time: MW 9:30-10:50

Lecture location: TBD

Section time: Mon 1-1:50

Section location: TBD

Instructors: Michael “Doc” Edge (he/him) and Peter Calabrese (he/him)

Instructor email: edgem@usc.edu, petercal@usc.edu

Edge office hours: TBD

Calabrese office hours: TBD

Teaching Assistant: TBD

TA email:

TA office hour:

Meeting ID:

Welcome! We are looking forward to working with you this semester.

Course Description

This is an upper-division course designed to introduce biologists to statistical theory for data analysis. Students will also learn basic programming skills in the statistical programming language R. The course is more mathematically demanding and more focused on general theory than BISC 305. At the same time, it is gentler and more targeted at biological data than courses that cover similar material in the math department, such as MATH 407 and 408. We will spend approximately 2/3 of the semester exploring simple linear regression, taking time to learn some statistical theory, to view linear regression from non-parametric/semi-parametric, likelihood-based, and Bayesian perspectives, and to implement methods in R. The remaining 1/3 of the semester will be a tour of some important techniques useful for describing, visualizing, and modeling different types of data, including from studies with multiple independent variables or dichotomous outcomes.

Textbook

We will be using *Statistical Thinking from Scratch: A Primer for Scientists*, by M.D. Edge, Oxford University Press, 2019. You can access an electronic copy of the book via the USC library. If you prefer to have a physical copy, you can buy one from Oxford University Press or your favorite bookseller.

The book also has a github directory (<https://github.com/mdedge/stfs>) with supplementary material, including all the code used in the book and solutions to all exercises.

Course Notes

In this course, we will take the time to learn one statistical method deeply first, and then we will add breadth at the end. This involves some mathematics and computer programming. Some of you may not have had math classes for a while and may have little experience programming. That will make the course a bit harder, but it is still possible to succeed with hard work and a good attitude. The grading system (see below) is designed to reward effort.

We will flip the classroom for some of the material, meaning that you will be expected to watch a taped lecture before class, and we will spend the class time learning actively. Slides for both traditional and flipped lectures will be posted.

Learning Goals

By course's end, our aim is that you will be able to:

- Discuss the philosophy involved in typical statistical estimation and inference, in which models are posited as data-generating processes with unknown parameters.
- Read and understand mathematical descriptions of simple statistical models.
- Explain the assumptions involved in justifying various views of the least-squares line, including a minimal "exploratory data analysis" view and views arising from semiparametric, parametric, and Bayesian models.
- Understand probabilistic and statistical concepts including expectation, variance, covariance, correlation, the law of large numbers, the central limit theorem, bias, consistency, efficiency, confidence intervals, p values, power, bootstrapping, permutation tests, likelihood, prior distributions, and posterior distributions.
- Design legible and informative data displays.
- Learn new methods for data analysis, such as linear regression, ANOVA, generalized linear models including logistic regression, principal component analysis, and linear mixed models, identifying principled reasons for choosing analysis methods.
- Explore the properties of statistical procedures using simulation and probability calculations.
- Use R to analyze and plot data, as well as write code to implement basic versions of procedures like bootstrapping and permutation testing.

Prerequisites

There are no specific requirements to enroll. The main requirement is that you have an interest in learning about using mathematics and computation to support scientific claims with data. Beyond that, comfort with algebra is very helpful. Some familiarity with the ways in which statistical analyses are used in research is helpful—if the words "mean," "median," "mode," "scatterplot," "standard deviation," "t-test," "confidence interval," are at least vaguely familiar, you are covered on this dimension. We will use some basic calculus, and we will be programming. Courses in these areas will likely help you feel comfortable initially but are not required.

If you have taken MATH 407 and 408 or equivalent courses, then the material in this class would be repetitive for you, and you are urged to take a different course.

Grading Policy

Your final grade will be calculated on the basis of a weighted average, with the weights

40% Homework
5% Participation
20% Term Paper
20% Final Exam
15% Midterm (take-home)

We will ask you to affirm that you have followed the rules for each exam.

Participation

Most lectures will include an activity or short response that's worth a participation point, to be submitted on blackboard. In case you need to watch some lectures asynchronously, you can submit the activity on blackboard before the next lecture.

Homework

There will be approximately 10 homework assignments during the semester, due every 1-2 weeks. Doing the homework will be your most important method for learning the material. Homework will be graded on a 0-3 scale, where a 0 indicates that a homework is missing or less than 50% complete, a 2 counts for full credit and represents a good effort on all problems, though some results may be wrong; and a 3 represents an exceptional effort. All "2"s would give you a perfect homework score. Scores of "3" will not happen often and are considered bonus.

Homework will be submitted via Turnitin on blackboard. Assignments will typically be due at 11:59 pm on Wednesdays (not every Wednesday, however). Assignments that are up to one week late will receive half credit. You are encouraged to work collaboratively on the homework, but please write your own solutions. We will drop your lowest homework score.

Questions

There is a general forum for questions and discussions on blackboard (under Tools > Discussion Board). Please ask questions about course content and general logistics here—if you have a question, someone else in the class likely has the same question, and answering it publicly will benefit everyone. If you have a logistics question related to your personal circumstances, please email an instructor or TA.

Software

We will use R, a programming language designed for statistical computing. R is available free online from the R Project website, <https://www.r-project.org/>. We recommend you also use RStudio, an interactive development environment designed for use with R. (The instructors will be using it.) RStudio is also free. (Download RStudio Desktop from <https://www.rstudio.com/products/rstudio/>.) RStudio requires an active R installation.

Course Schedule (Subject to change)

Jan 10: (Edge): Intro, course policies. The regression line as a motivating problem.
Reading: *STFS* Prelude and Chapter 1.

Unit 1: Mathematical and computational tools for statistics

Jan 12: (Calabrese/Edge): The importance of data quality; the least-squares line
Reading: *STFS* chapter 3. **Note that we are reading chapter 3 before chapter 2.**

Week 2

Jan 17: Martin Luther King, Jr. Day, no class

Jan 19: (Edge/Calabrese): The statistical programming language R, part 1: interacting with R and exploratory data analysis. Note: lectures will be “flipped” Jan 19-26. Please watch didactic material before class; we will work on programming in groups during class time.

Reading: *STFS* chapter 2.

Week 3

Jan 24 (Edge/Calabrese): R, part 2 - Functions and Loops.

Reading: *STFS* Appendix B.

Jan 26 (Edge/Calabrese): R, part 3 - data input/output; R markdown and RStudio notebooks.

Reading: *STFS* Appendix B.

Week 4

Jan 31 (Calabrese): Probability 1. (Foundations, Axioms, independence, conditional probability, Bayes’ Theorem)

Reading: *STFS* chapter 4 (through the end of section 4.2 / Box 4-2, pp 38-48)

Feb 2 (Calabrese): Probability 2. (Discrete and continuous random variables, pdfs and cdfs, distribution families)

Reading: *STFS* chapter 4 (sections 4.4-4.8, pp 48-58)

Week 5

Feb 7 (Calabrese): Probability 3. (Expectation, Variance, and the law of large numbers)

Reading: *STFS* chapter 5 (through the end of section 5.2, pp 60-68)

Feb 9 (Calabrese): Probability 4. (Correlation and covariance; The central limit theorem)

Reading: *STFS* chapter 5 (sections 5.3 and 5.5)

Week 6

Feb 14 (Calabrese): Probability 5. (conditional distributions; a model for linear regression)

Reading: *STFS* chapter 5 (sections 5.4, 5.6-5.7)

Unit 2: Basic statistical theory

Feb 16 (Edge): Properties of Estimators: Bias, Variance, Mean Squared Error, and Consistency.

Reading: *STFS* interlude; chapter 6 through the end of section 6.4.

Week 7

Feb 21: President's Day, no class

Feb 23: (Edge): Properties of Estimators: Efficiency and Robustness. Decision Theory.

Reading : *STFS* chapter 6, sections 6.5-6.10.

Week 8

Feb 28 (Edge): Standard error and confidence intervals.

Reading: *STFS* chapter 7 through the end of section 7.2.

March 2 (Edge): p values and hypothesis tests

Reading: *STFS* chapter 7, sections 7.2-7.4

Week 9

March 7 (Calabrese): Multiple testing in genomics + midterm Q&A

Reading: N/A

Take-home midterm released after section; due by 11:59 pm on Wednesday, March 9.

March 9 (Edge): Power and effect size. Criticisms of NHST.

Reading: *STFS* chapter 7, sections 7.6-7.9 (skip optional section 7.5)

Take-home midterm due 11:59pm

March 14 & 16: Spring Break, no class

Unit 3: Three major approaches to estimation and inference

Week 10

March 21 (Edge): Plug-in estimators, the method of moments, and the bootstrap.

Reading: *STFS* chapter 8 through the end of section 8.2.

March 23 (Edge): Permutation tests.

Reading: *STFS* chapter 8, sections 8.3-8.5.

Week 11

March 28 (Edge): Maximum-likelihood estimation.

Reading: *STFS* chapter 9, through section 9.2; skip optional section 9.2.2.

March 30 (Edge): Wald test, score test, and likelihood-ratio test.

Reading: *STFS* chapter 9, sections 9.3-9.5.

Week 12

April 4 (Edge): The Bayesian Alternative: Priors and posteriors.

Reading: *STFS* chapter 10.

Unit 4: Models for data analysis

April 6 (Edge): Assessing linear regression assumptions, multiple linear regression

Reading: *STFS* Postlude, through the end of section Post.2.1

Week 13

April 13 (Edge): Multiple regression / Generalized linear models

Reading: *STFS* sections Post.2.2-end

April 15 (Calabrese): Special cases and relatives of linear regression 1

Reading: Course notes (will be posted on blackboard)

Week 14

April 20 (Calabrese): Special cases and relatives of linear regression 2

Reading: Course notes posted on blackboard

April 22 (Edge): Causal inference

Reading: None

Week 15

April 27 (Edge): Statistics in Society: Eugenics as a cautionary tale

Reading: None

April 29 (Calabrese): Neural networks

Reading: None

Term paper due April 29th by 5 pm.

Our final exam slot is Friday, May 6th, 8-10am

Term paper guidelines

For most of you, the term paper will involve modeling and analyzing a data set of your choice. You can consider a public dataset, such as one included in an R package or posted on the R directory “Free data sets” page: <https://r-dir.com/reference/datasets.html>. If you are having trouble finding a suitable data set, please check in with an instructor. Alternatively, if you are conducting research, you may have your own dataset that you would like to analyze. Your analysis must be original; you cannot repeat an existing analysis of a public dataset.

A second option for the term paper is to run a simulation study, wherein you study the properties of some statistical procedure(s) when applied to hundreds or thousands of simulated datasets with known properties. You could use such a study to compare the properties of different statistical procedures, such as bootstrap-based vs. normal theory confidence intervals or linear mixed models vs. generalized estimating equations. If you are interested in completing a simulation study rather than an analysis of existing data, please check in with an instructor, and we will help you make a plan.

The below outline describes a term paper involving a data analysis. Compared with a standard research paper in biology, your term paper will have abbreviated introduction, methods, and discussion, but it will have substantial detail in the results.

Term Paper Outline. (The approximate rubric below is out of 90. There will also be 10 points awarded for mechanics and style. If there are mechanics/style problems that make any of the below sections difficult to comprehend in certain places, then the relevant section scores may suffer as well.)

Introduction. State the question/hypothesis that motivates your examination of the dataset, with a brief theoretical framework. Finish the introduction with a statement that says how this data set addresses the question. This introduction should be brief—something like 1-3 paragraphs—and include

2-5 references. The goal here is *not* to provide a thorough literature review but just to tell the reader why you're analyzing this dataset and what question you'll be trying to answer (10 points).

Methods. Describe how the data were collected. One short paragraph. The idea is to highlight the key features of the data that will help the reader interpret the results, not necessarily to provide all the details necessary for replicating the study.

Exploratory data analysis. Create figures and tables to illustrate the major patterns in the data. The figures should have appropriate axis labels and, if applicable, either a legend or labels directly on the plot. Each figure will also need a caption. We expect 2 to 5 figures. Embed the figures in your text, commenting on the main things you can learn from the summaries. (10 points, including the R code, which will be submitted as an Appendix.)

Model Description. State what your competing models are both in English and with equations. The competing models can be just the null model and a particular model that corresponds to your hypothesis as in a classical *hypothesis testing* procedure. But you can also decide that it is more relevant to compare different complex models. (10 points).

Model estimation and analysis. We will cover three major perspectives for data analysis this quarter: nonparametric/semiparametric, parametric frequentist (maximum likelihood), and Bayesian. Unless you get a special exception for your project, you will need to analyze your data using techniques from *two* of these three families of methods. That is, you will complete two parallel sets of analysis in two different frameworks.

For example, if you use linear regression, you might estimate the parameters using least squares and then compute standard errors *both* using bootstrapping (non/semiparametric) and using typical normal-theory standard errors (parametric frequentist), then continue to build confidence intervals and compute p values using each set of standard error estimates. Or perhaps you are using a linear mixed model, and you compute point and interval estimates using maximum likelihood (maximum likelihood) and also by maximum a posteriori methods (Bayesian), basing prior distributions for the parameters on past work. Many other combinations are possible. Whatever you do, give and interpret complete results using both analysis methods, and compare and contrast the relative advantages and disadvantages of the analysis approaches you use. (35 points)

Analysis of model appropriateness. Assess the appropriateness of the model(s) you used. Do your data seem to violate the model's main assumptions? Use data displays and/or hypothesis tests to make your case, and discuss the implications. (10 points)

Conclusion. Write a brief conclusion of your modeling efforts, the statistical analysis and the implications for the question that was raised in the introduction. One or two paragraphs. (10 points)

Appendix. R code.

Statement on Academic Conduct and Support Systems

Academic Conduct:

Plagiarism – presenting someone else's ideas as your own, either verbatim or recast in your own words – is a serious academic offense with serious consequences. Please familiarize yourself with the discussion of plagiarism in SCampus in Part B, Section 11, "Behavior Violating University Standards" policy.usc.edu/scampus-part-b. Other forms of academic dishonesty are equally unacceptable. See additional information in SCampus and university policies on scientific misconduct, policy.usc.edu/scientific-misconduct.

Support Systems:

Counseling and Mental Health - (213) 740-9355 – 24/7 on call
studenthealth.usc.edu/counseling

Free and confidential mental health treatment for students, including short-term psychotherapy, group counseling, stress fitness workshops, and crisis intervention.

National Suicide Prevention Lifeline - 1 (800) 273-8255 – 24/7 on call
suicidepreventionlifeline.org

Free and confidential emotional support to people in suicidal crisis or emotional distress 24 hours a day, 7 days a week.

Relationship and Sexual Violence Prevention and Services (RSVP) - (213) 740-9355(WELL), press “0” after hours – 24/7 on call
studenthealth.usc.edu/sexual-assault

Free and confidential therapy services, workshops, and training for situations related to gender-based harm.

Office of Equity and Diversity (OED)- (213) 740-5086 | Title IX – (213) 821-8298
equity.usc.edu, titleix.usc.edu

Information about how to get help or help someone affected by harassment or discrimination, rights of protected classes, reporting options, and additional resources for students, faculty, staff, visitors, and applicants. The university prohibits discrimination or harassment based on the following *protected characteristics*: race, color, national origin, ancestry, religion, sex, gender, gender identity, gender expression, sexual orientation, age, physical disability, medical condition, mental disability, marital status, pregnancy, veteran status, genetic information, and any other characteristic which may be specified in applicable laws and governmental regulations. The university also prohibits sexual assault, non-consensual sexual contact, sexual misconduct, intimate partner violence, stalking, malicious dissuasion, retaliation, and violation of interim measures.

Reporting Incidents of Bias or Harassment - (213) 740-5086 or (213) 821-8298
usc-advocate.symplicity.com/care_report

Avenue to report incidents of bias, hate crimes, and microaggressions to the Office of Equity and Diversity | Title IX for appropriate investigation, supportive measures, and response.

The Office of Disability Services and Programs - (213) 740-0776
dsp.usc.edu

Support and accommodations for students with disabilities. Services include assistance in providing readers/notetakers/interpreters, special accommodations for test taking needs, assistance with architectural barriers, assistive technology, and support for individual needs.

USC Support and Advocacy - (213) 821-4710
uscса.usc.edu

Assists students and families in resolving complex personal, financial, and academic issues adversely affecting their success as a student.

Diversity at USC - (213) 740-2101
diversity.usc.edu

Information on events, programs and training, the Provost's Diversity and Inclusion Council, Diversity Liaisons for each academic school, chronology, participation, and various resources for students.

USC Emergency - UPC: (213) 740-4321, HSC: (323) 442-1000 – 24/7 on call

dps.usc.edu, emergency.usc.edu

Emergency assistance and avenue to report a crime. Latest updates regarding safety, including ways in which instruction will be continued if an officially declared emergency makes travel to campus infeasible.

USC Department of Public Safety - UPC: (213) 740-6000, HSC: (323) 442-120 – 24/7 on call

dps.usc.edu

Non-emergency assistance or information.