

# Statistics 694: Research in Statistics and Biostatistics (2 units)

Department of Statistics and Biostatistics, CSU East Bay

Prof. Eric A. Suess

2023-08-25

**Course Description:** A collaborative research experience on a research topic designated by the instructor in the field of Statistics. Students conduct a literature search, develop a research proposal, complete a research project leading to a professional presentation, paper, or report.

**Lecture:** F noon to 2pm, South Science 302

**Instructor:** Prof. Eric A. Suess

**Office:** NSc 319 **Phone:** 510-885-3879 **e-mail:** eric.suess@csueastbay.edu

**Office Hours:**

- Thursday online 10:00 - 11:00am Here is the Zoom link: <https://csueb.zoom.us/j/89647106693>
- Friday NSc 319 2:00 - 3:00pm
- or by appointment

**Communicating:**

Email is the preferred method of communication. Class website will be updated weekly with class topics, homework assignment, and other useful information. Assignment grades will be provided in Canvas.

**Class Website:** [cox.csueastbay.edu/~esuess/statistics694](http://cox.csueastbay.edu/~esuess/statistics694)

**Course Materials:** Access to a modern computer and permission to install software software. Access to the internet.

**Required Texts:**

- Wickham, Golemund [R for Data Science, 2nd edition](#)
- Wickham, [Advanced R](#)

- Ismay, Kim, [ModernDive](#)
- Phillips, [Yarr](#)
- Speegle, [Probability, Statistics, and Data: A Fresh Approach Using R](#)
- Boehmke, Greenwell, [Hands-On Machine Learning with R](#)
- Kross, [Unix Workbench](#)
- Silge, Robinson, [TidyText](#)
- Wickham, [Mastering Shiny](#)
- Sievert, [Interactive web-based data visualization with R, plotly, and shiny](#)
- Janssens, [Data Science at the Command Line](#)

**Reference Texts:**

- Baumer, Kaplan, Horton, [Modern Data Science with R, 3rd edition](#), CRC Press, 2023.

**Further References:**

- Bryan, Hester, [Happy Git and GitHub for the useR](#)
- VanderPlas, [Python Data Science Handbook](#)
- Klok, Nazarathy, [Statistics with Julia](#)
- Garrels, [Bash Beginners Guide](#)

**Material To Be Covered:**

This is the **zero** course in the sequence of Data Science courses offered by the Department of Statistics and Biostatistics for the MS Data Science Concentration. The Data Science courses are specifically for registered students in the MS Statistics program.

The sequence of courses are:

0. Stat. 694 Research in Statistics and Biostatistics - alternatively: Innovation, Data Technologies, Data Science Workgroup, Foundations in R for Data Science
1. Stat. 650 Advanced R for Data Science
2. Stat. 651 Data Visualization
3. Stat. 652 Statistical Learning
4. Stat. 653 Statistical Natural Language Processing
5. Stat. 654 Introduction to Applied Deep Learning

These courses are intended to be taken in order as they build upon each other, but you can discuss taking the courses out of order with instructor approval.

(Alternatively, students who have completed the 650 sequence, this course can be used to reinforce your R Notebook skills, your understanding of data structures in R, expand your experience with the datasets used in the 650 sequence, and to research topics beyond the classes. Start to learn Python.)

The topics in this course will be drawn from a variety of on-line books, blog posts, YouTube videos, and Github. For each class there will be other supporting references. The intention

is to develop a collection of resources for the 650 sequence. And to have experience working with the main datasets that will be part of the 650 sequence, to try some and become aware of many modern data technologies that are commonly used.

The main topics for this class may include:

- R Notebooks, the file system, your computer system, benchmarking, parallel processing, cluster computing, cloud computing, super computers, data.frames, data.table, lists, Tidyverse, data, structured data, compressed data, zip and unzip. Python Notebook: Google CoLab.
- Structured vs unstructured data, AutoEDA, AutoML, h2O, Tensorflow/Keras, SQL, noSQL, SQLite, MongoDB, Small Data, Medium Data, Big Data, Spark.
- Unstructured data, images, music, video, APIs, Github, Slack, Loom.
- R, Python, Julia, Bash, SQL Jupyter Notebooks, Anaconda, Kaggle.
- Unstructured data, text, e-books, many small files into one dataset, break one Big file into smaller files, arrow.
- Other Topics: Shiny, R packages, using R to set up an API.
- Cloud Computing: Linux, [Data Linux](#), BSD, Docker, Amazon AWS, Microsoft Azure, Google GCP, IBM Watson, Digital Ocean, Edge Computing, IoT, RaspberryPi, FitBit, iPhones, Android Phones.
- AI: Deep Learning, Transformers, Stable Diffusion, LLMs, GTP, Bing Chat, Hugging Face, [RTutor](#), responsible use of AI for learning and doing Data Science.

The datasets for this class may include:

- iris, palmerpenguins, mtcars, diamonds, gapmider, nycflights, NHANES, MNIST, Boston Housing, cat-and-dogs
- bikeshare data, Lyft Bay Wheels, NYC Taxi data, Airline On-Time Statistics, Chicago Crime data, Fannie Mae
- Twitter, Gutenberg, NASA, Genius, Spotify, Arxiv, Internet Archive.

### **Technical Requirements:**

Access to a modern computer and permission to install software, R and RStudio. Access to the internet.

### **Homework:**

Homework will be assigned weekly. Homework will be “due” on the following Fridays, which means you should complete the homework and come to class prepared to ask questions. Homework will be “collected” through Canvas.

### **Project:**

Idea, Plan & Development, Data Product, Presentation/Video.

**Grading:**  $\geq 90\%$  A,  $\geq 80\%$  B,  $\geq 70\%$  C,  $\geq 60\%$  D,  $<60\%$  F

- Homework 50%
- Project 50%

**Policy on Make-up Exams:** You are expected to take the quizzes and exams at the scheduled times. In case of genuine emergency, illness or hardship, for which you can present written documentation I may agree to arrange for a make-up exam. Make-up exams must always be arranged BEFORE the regular exam is given and always take place AFTER the regular exam. Quizzes may not be made up!

### **Statistics 650 SLOs**

#### **Student Learning Outcomes (SLO's):**

Students graduating with an M.S. in Statistics from Cal State East Bay will be able to:

1. Apply statistical methodologies, including a) descriptive statistics and graphical displays, b) probability models for uncertainty, stochastic processes, and distribution theory, c) hypothesis testing and confidence intervals, d) ANOVA and regression models (including linear, and multiple linear) and analysis of residuals from models and trends at the Master's level.
2. Derive basic theory underlying these methodologies.
3. Model practical problems for solutions using these methodologies.
4. Produce relevant computer output using standard statistical software and interpret the results appropriately.
5. Communicate statistical concepts and analytical results clearly and appropriately to others; and,
6. Employ theory, concepts, and terminology at a level that supports lifelong learning of related methodologies.