

# Stats 598z: Midterm exam 2

---

**Important:**

Write your name and PUID on all sheets, and include the number of sheets

There are 6 questions, questions 1 and 6 have 9 points, the rest 8 points

Attempt all questions, and when appropriate include a brief justification of your answer

Don't spend time polishing your answers as the main idea is more important

---

- A dataframe `my_class` has columns `puid`, `degree`, `major`, and `grade`. In the following, use commands from `tidyverse` for full points, but you will get most points if you just use base R.
    - Write code to sort the rows of the dataframe in decreasing order of `grade`.
    - Write code to calculate the number of students of every (`degree`, `major`) pair as well as the mean score.
    - Write code to calculate the number of students of every `major` with a grade greater than 90.
    - Write code to plot histograms of `grade`, with each (`degree`, `major`) pair having its own color.
  - Provide short examples displaying some data in both tall and wide format.
- You have a vector of filenames `my_files`. Write down the R regular expression command to find:
  - files whose names contain “`data`” in them
  - files whose names end with “`.txt`”
  - files whose names contain two or more numbers in a row.
  - Suppose all filenames contain a date, taking the form `MM/DD/YYYY` (and no other numbers). Write R code to replace the date in the names to `DD/MM/YYYY`.
- Briefly describe the simplex (or Nelder-Mead) algorithm at a high-level.
  - In the context of object-oriented programming, describe 1) inheritance 2) polymorphism 3) generic functions.
  - What are the outputs for `2 %>% '/(4)` and `2 %>% '/(4,.)`?
- Explain briefly what reactive programming is.
  - For the dataframe `my_class`, use `ggvis` to plot a histogram of grades, setting the number of bins with a slider.
  - Explain what the command `log1p(x)` does, and how it is different from `log(1+x)`.
  - Explain briefly why `.3*4 == 1.2` might behave unexpectedly, and how you might do this correctly.
- What is the motivation for LASSO? Write down the LASSO loss-function.
  - What is coordinate descent?
  - Recall that LASSO tries to find the `w` that minimizes the LASSO loss-function. Suppose you wanted to find the average value of the LASSO loss function for a fixed `X`, `y`, `lambda`, with `w` distributed as a Gaussian random variable. How would you do this? Provide a few lines of R code.
- Explain what the `set.seed()` function does, and why it is useful.
  - Consider rolling a fair die until the number 6 shows up. The number of rolls required is a random number. Write R code to get a Monte Carlo estimate of the average the value of this quantity.
  - Describe importance sampling at a high-level. Explain (e.g. with an example) how it is useful for estimating quantities involving rare events.