

## Stats 598z: Midterm exam 2

---

### Important:

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

There are 7 questions, each for 5 points (but not all equally easy)

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

---

- `my_vec2` is a vector of integers. Write a few lines of R to return the number of times two consecutive numbers occur in successive positions. E.g. `(5, 6, 7, 2, 4, 5)` returns 3. Use vectorization for full points.
  - `my_vec1` is a vector consisting of a sequence of numbers that are interpreted as `(len1, obs1, obs2, ..., obs1en1, len2, obs1, obs2, ..., obs1en2, ...)`. Thus, the first number gives the length of a set of observations that follows it, these in turn are followed by the length of the next set, and so on. Write a few lines of R code to return the mean of each set, so that `(3, 1, 2, 3, 2, 5, 3)` returns `(2, 4)`. You can use loops.
- You have a vector `my_words`. Write down the R regular expression command to find components
  - containing `m` followed by one or more `c`'s.
  - starting with a vowel.
  - starting and ending with a vowel.
  - containing `".txt"`.
  - Write R code to replace all `".txt"` with `".doc"`.
- What is coordinate descent? What is gradient descent? Briefly explain possible advantages and disadvantages.
  - What is the role of  $\lambda$  in LASSO. Briefly explain how you will use cross-validation to choose  $\lambda$ .
  - For a fixed dataset, you want to solve a series of LASSO problems for  $\lambda = 0, .1, .2, .5, 1, 2, 5$  and 10. Explain briefly how you can order these problems for efficiency.
- Briefly describe some of the advantages of object-oriented programming.
  - You have `my_list`, a list of vectors of numbers. What does `my_list[1]` return (and what type of data object is this)? How would you use object-oriented programming so that `my_list[i]` returns the mean of the `i`th vector of `my_list`? For full points, provide necessary R code.
- What is tidy data? What are tall and wide formats of dataframes? Give a simple example.
  - A data frame `population` has columns `name`, `city`, `age`, `occupation`, `gender` and `income`, each row giving these quantities for a different individual. Give R code to convert it into a new dataframe giving the average income for each occupation and for each gender (averaged across all people, cities, and ages). What if you wanted to do this for each occupation, gender and age decade as well (i.e. 0-9, 10-19 etc.)? To get full points, use `melt()`, `dcast()`, `*ply` or vectorization, but if not sure, use loops.
- Describe importance sampling at a high-level. Explain (e.g. with an example) how it is useful for estimating quantities involving rare events.
  - Briefly describe any Monte Carlo method to estimate  $\int_a^b \log(x) dx$ . Provide a few lines of R code.
- Explain why Metropolis-Hastings might be better than importance/rejection sampling.
  - Given a 100 coins, you want to calculate the average number of heads given that the number of heads is greater than 70:  $\mathbb{E}[\#H | \#H > 70]$ . Briefly describe a Metropolis-Hastings algorithm to do this, giving the initialization, the proposal distribution and the acceptance probability.