

# CSCI 1051 Homework 4

January 29, 2024

## Submission Instructions

Please upload your solutions by **5pm Friday February, 2024**.

- You are encouraged to discuss ideas and work with your classmates. However, you **must acknowledge** your collaborators at the top of each solution on which you collaborated with others and you **must write** your solutions and code independently.
- Your solutions to theory questions must be typeset in LaTeX or markdown. I strongly recommend uploading the source LaTeX (found [here](#)) to Overleaf for editing.
- I recommend that you write your solutions to coding questions in a Jupyter notebook using Google Colab.
- You should submit your solutions as a **single PDF** via the assignment on Gradescope. You can enroll in the class using the code GPXX7N.
- Once you uploaded your solution, **mark where you answered each part of each question**.

## Problem 1

Consider the linear regression problem with  $n \geq d$ . Let  $\mathbf{A} \in \mathbb{R}^{n \times d}$  be a feature matrix and  $\mathbf{b} \in \mathbb{R}^n$  be a target vector. The regression problem is to find a minimizing vector

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{Ax} - \mathbf{b}\|_2^2.$$

You previously showed that the optimal solution is  $\mathbf{x}^* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$ . In this problem, you will compare computing the optimal solution exactly to computing it approximately using the fast Johnson-Lindenstrauss transform. We will use the MNIST dataset to build  $\mathbf{A}$  and  $\mathbf{b}$ . The MNIST dataset consists of  $28 \times 28$  pixel handrawn digits of numbers with the corresponding label.

### Part 1 (1 point)

Using the code I provide in `regression.py`, compute the exact solution  $\mathbf{x}^*$  and the mean squared error

$$\frac{1}{n} \|\mathbf{Ax} - \mathbf{b}\|_2^2.$$

If your code is anything like mine, it will be slow and return a a terrible solution due to *round off error*.

### Part 2 (2 points)

Now implement the fast JL transform as described in class. In particular, compute  $\mathbf{\Pi A} = \mathbf{SHDA}$  one column of  $\mathbf{A}$  at a time. Recall that  $\mathbf{S}$  is a sampling matrix,  $\mathbf{H}$  is a Hadamard, and  $\mathbf{D}$  is a diagonal matrix with a random sign.

When you are done, compute the mean squared error of your solution and comment on how it compares to the “exact” solution.

**Hint:** Computing  $\mathbf{H}$  is too expensive so write a function to compute  $\mathbf{HDx}$  using recursion. You can speed up the recursion by checking if there are any non-zeros in the vector.

## Problem 2

Thank you for taking this class with me! As I've mentioned, I love randomized algorithms for data science because the topic combines beautiful math with *interesting* applications. I know I have a lot to improve and I would love your feedback on what went well and what could have gone better! Here are some of the aspects of the course I've thought a lot about but you can give me feedback on anything.

- **Content:** What topics did you like? What would you like to have covered? What would you be okay skipping?
- **Difficulty:** How was the difficulty of the class?
- **Daily Check in Forms:** What do you think about the daily check in forms?
- **Group Activities:** What did you think about the group activities?
- **Content Review:** What did you think about the content review the next day?
- **Accessibility:** How accessible was I as a teacher? Did you feel comfortable asking me questions? Did I give enough or too many hints when asked about problems?
- **Afternoon Problem Solving:** What did you think about the afternoon problem solving session?
- **Self-Grade and LaTeX:** What did you think about the self-grade and writing your solutions in LaTeX?
- **Typed Notes and Slides:** What did you think about having the typed notes available online? How about the slides?

### Part 1 (1.5 points)

Please tell me what you liked about the class so I can do more of it in the future.

### Part 2 (1.5 points)

Please tell me what I could improve to make the experience better.