

## CS 475/675 Machine Learning: Homework 1

## Supervised Learning 1

Due: Thursday, September 24, 2020, 11:59 pm US/Eastern

65 Points Total

Version 1.0

**Make sure to read from start to finish before beginning the assignment.**

## 1 Homeworks

This semester we are introducing a new homework system for the course. This new system is based on feedback from previous semesters and adaptations given that the course is now fully online. We anticipate making changes over the course of the semester based on your feedback.

We will have two types of homework assignments this semester: Written and Programming. We expect to alternate the type of assignment, with a total of 8 assignments being worth a total of 400 points (out of 1000 available points). In general, homeworks will be due on Thursday evening at 11:59pm, and homeworks will be released 10-14 days before they are due. There will be an assignment due most weeks.

Our current plan, subject to change, is that the written assignments will be worth 65 points and the programming worth 35 points. We will begin (i.e. this assignment) with a written assignment. Each assignment will contain a version number at the top. While we try to ensure every homework is perfect when we release it, errors do happen. When we correct these, we'll update the version number, post a new PDF and announce the change. Each homework starts at version 1.0.

### 1.1 Written Assignment

The written assignment will contain two parts:

1. **Analytical:** These analytical questions will consider topics from the course. These will include mathematical derivations and analyses. Your answers will be entirely based on written work, i.e. no programming.
2. **Lab:** In the lab portion of the assignment, you will apply machine learning concepts to gain experience working with data from different domains. Labs will typically involve a Python notebook, a write up and an applied exploration of topics covered in the class.

The point total for each portion of the homework will be listed in the assignment. Written assignments will be submitted as PDFs. See below for more details about what to submit.

## 1.2 Collaboration Policy

The course policy is that, unless otherwise specified, all work must be your own. See the about page on the course website for more details.

**For this assignment, we strongly recommend you work with a partner.** You and your partner will make one submission for the two of you on Gradescope (make sure to include your partner when you submit). You and your partner will receive the same grade, so please choose your partner carefully.

You can only work in teams of one or two (not more). Your partner can be anyone from either section (01 or 02) or course (475 or 675) provided that both of you are taking the course for credit (not audit). We *highly* recommend that you do every part of the assignment together instead of splitting it up. You can work on the same Overleaf document and think through the questions together. You probably want to work with the same partner for the semester (*only* for assignments where collaboration is allowed) but it is not a requirement.

## 1.3 What to Submit

For this assignment you will submit the following.

1. **Analytical.** You will submit your analytical solutions to Gradescope. **Your writeup must be compiled from L<sup>A</sup>T<sub>E</sub>X and uploaded as a PDF.** The writeup should contain all of the answers to the analytical questions asked in the assignment. Make sure to include your name in the writeup PDF and to use the provided L<sup>A</sup>T<sub>E</sub>X template for your answers following the distributed template. You will submit this to the assignment called “Homework 1: Supervised Learning 1: Analytical”.
2. **Lab Python Notebook** You will submit your Python notebook as a PDF by going to File → Export via PDF or File → Export via PDF via LaTeX. Once you download the pdf, open the file to ensure that the plots show up. You will submit this to the assignment called “Homework 1: Supervised Learning 1: Lab”.
3. **Lab Data** You will submit your data and associated files as a zip file. You will submit this to the assignment called “Homework 1: Supervised Learning 1: Lab Data”.

You will need to create an account on gradescope.com and signup for this class. The course is <https://www.gradescope.com/courses/153788>. Use entry code M83BRX. **You must either use the email account associated with your JHED, or specify your JHED as your student ID.** See this video for instructions on how to upload a homework assignment: [https://www.youtube.com/watch?v=KMPoby5g\\_nE](https://www.youtube.com/watch?v=KMPoby5g_nE).

## 1.4 Questions?

Remember to submit questions about the assignment to the appropriate group on Piazza: [piazza.com/jhu/fall2020/601475](https://piazza.com/jhu/fall2020/601475).

## 2 Analytical (30 points)

Please see the accompanying `homework1_template.tex` file for the analytical questions for this assignment. There is space provided in that file for you to type your answers

in  $\text{\LaTeX}$  after each question. **Do not edit the file in any way except to add your answers.** Gradescope assumes that the PDF will exactly match our template except for your solutions.

In addition to completing the analytical questions, your assignment for this homework is to learn  $\text{\LaTeX}$ . All homework writeups must be PDFs compiled from  $\text{\LaTeX}$ . Why learn  $\text{\LaTeX}$ ?

1. It is incredibly useful for writing mathematical expressions.
2. It makes references simple.
3. Many academic papers are written in  $\text{\LaTeX}$ .

The list goes on. Additionally, it makes your assignments much easier to read than if they are written by hand or if you complete them in Word.

We realize learning  $\text{\LaTeX}$  can be daunting. Fear not. There are many tutorials on the Web to help you learn. We recommend using `pdflatex`. It's available for nearly every operating system. As the semester progresses, you'll no doubt become more familiar with  $\text{\LaTeX}$ , and even begin to appreciate using it.

Be sure to check out this cool  $\text{\LaTeX}$  tool for finding symbols. It uses machine learning! <http://detexify.kirelabs.org/classify.html>

For each homework analytical we will provide you with a  $\text{\LaTeX}$  template. You **must use the template**. The template contains detailed directions about how to use it.

Please open the template to view the analytical questions.

### 3 Lab (35 points)

In this assignment you will be creating a dataset for supervised learning. You will also become familiar with some popular off-the-shelf machine learning tools people use in practice. You will evaluate your dataset using the framework we introduced in class:

1. Is there a well-defined problem?
2. Does an easy solution exist for the problem?
3. Do you have large amounts of high quality data?
4. Can you meaningfully evaluate results? What would the loss function measure?
5. Is using machine learning for this problem justified?

You should create a dataset for a problem where applying machine learning is challenging, *but still possible*. In other words, a supervised machine learning algorithm should be able to generalize from a training set of  $(x, y)$  pairs to make predictions for unseen  $x$  examples.

Make sure to think through the ethical implications of the data you are collecting<sup>1</sup>. Beyond this course, as future researchers and practitioners of machine learning, you must consider ethical implications of your work. We'll learn more about this over the semester.

---

<sup>1</sup>Not sure how to think through these ethical implications? Start by reading this Medium article: <https://medium.com/@mrtz/how-big-data-is-unfair-9aa544d739de>

### 3.1 Identifying the data source

You are free to use data from any domain of interest. We provide some examples of sources of data.

- If you are interested in text data, Wikipedia is a great starting place (<https://meta.wikimedia.org/wiki/Datasets>). Let's consider a Wikipedia document as our example  $x$ . Then, we may be interested in predicting  $y$ , where  $y$  is the number of page revisions, the number of authors, the number of page views, the topic of the page, the language the page is written in, etc.
- If you're interested in image data, consider exploring this collection of open image datasets for inspiration: <https://blogs.ntu.edu.sg/openimagecollections/browse/#collections>. Let's consider an image as our example  $x$ . Then, we may be interested in predicting  $y$ , where  $y$  is the year the image was created, the artist who created the image, the medium of the image, etc.
- If you are interested in public policy, consider exploring datasets produced by the US government (<https://www.data.gov/>) and by the Baltimore City government (<https://data.baltimorecity.gov/>). There are a number of directions to take that address social problems.
- If you are interested in health, consider exploring datasets produced by the CDC ([https://www.cdc.gov/nchs/data\\_access/ftp\\_data.htm](https://www.cdc.gov/nchs/data_access/ftp_data.htm)) or related to COVID-19 (<http://www.socialmediaforpublichealth.org/covid-19/resources/>).
- Also see this repository of structured data (<https://docs.google.com/spreadsheets/d/1wZhPLMCHKJvw0kP4juclhjFgqIY8fQFMemwKL2c64vk>). There are many domains and data formats represented.

You may use data from a combination of sources you identify, and you should have a clear idea of the problem you are trying to solve with the data you are collecting.

**You may not just download an existing dataset to use!** We want you to create a new dataset for supervised learning. Pick something of interest to you, and that others in the class would find interesting. The labels in your data can be automatically derived from the online source (e.g. how many Wikipedia page views?) or manually applied by you (e.g. did I like this song on Spotify?). If you find some features you like for  $x$ , you should at least find some other labels to predict for  $y$ .

### 3.2 Creating the dataset

Open the Jupyter notebook `homework1_lab.ipynb`. This notebook will walk you through defining your problem, creating the dataset, exploring your data, and running some basic machine learning algorithms. There are questions that should be answered inline within the notebook.

You will hand in both the Python notebook, which contains answers to the questions, and the dataset you create.