

---

# CS5785 Homework 1

---

The homework is generally split into programming exercises and written exercises.

This homework is due on **September 26, 2018 at 11:59 PM PST**. Upload your homework to [CMS](#), unless you can't upload via cmsx as discussed previously with ta's - in that case send a slack message with the hw. Please upload the submission as two files - a complete submission should include:

1. A write-up as a single .pdf file
2. Source code for all of your experiments (AND figures) zipped into a single .zip file, in .py files if you use Python or .ipynb files if you use the IPython Notebook. If you use some other language, include all build scripts necessary to build and run your project along with instructions on how to compile and run your code.

The write-up should contain a general summary of what you did, how well your solution works, any insights you found, etc. On the cover page, include the class name, homework number, and team member names. You are responsible for submitting clear, organized answers to the questions. You could use online  $\LaTeX$  templates from [Overleaf](#), under "Homework Assignment" and "Project / Lab Report".

Please include all relevant information for a question, including text response, equations, figures, graphs, output, etc. If you include graphs, be sure to include the source code that generated them. Please pay attention to Slack for relevant information regarding updates, tips, and policy changes. You are encouraged (but not required) to work in groups of 2.

## IF YOU NEED HELP

There are several strategies available to you.

- If you ever get stuck, the best way is to ask on Slack. That way, your solutions will be available to the other students in the class.
- Your instructor and TAs will offer office hours<sup>1</sup>, which are a great way to get some one-on-one help.
- You are allowed to use well known libraries such as `scikit-learn`, `scikit-image`, `numpy`, `scipy`, etc. in this assignment. Any reference or copy of public code repositories should be properly cited in your submission (examples include Github, Wikipedia, Blogs).

---

<sup>1</sup><https://cornelltech.github.io/cs5785-fall-2018/contact.html>

## PROGRAMMING EXERCISES

### 1. Digit Recognizer

- (a) Join the [Digit Recognizer](#) competition on Kaggle. Download the training and test data. The competition page describes how these files are formatted.
- (b) Write a function to display an MNIST digit. Display one of each digit.
- (c) Examine the prior probability of the classes in the training data. Is it uniform across the digits? Display a normalized histogram of digit counts. Is it even?
- (d) Pick one example of each digit from your training data. Then, for each sample digit, compute and show the best match (nearest neighbor) between your chosen sample and the rest of the training data. Use  $L_2$  distance between the two images' pixel values as the metric. This probably won't be perfect, so add an asterisk next to the erroneous examples (if any).
- (e) Consider the case of binary comparison between the digits 0 and 1. Ignoring all the other digits, compute the pairwise distances for all genuine matches and all impostor matches, again using the  $L_2$  norm. Plot histograms of the genuine and impostor distances on the same set of axes.
- (f) Generate an ROC curve from the above sets of distances. What is the equal error rate? What is the error rate of a classifier that simply guesses randomly?
- (g) Implement a K-NN classifier. (You cannot use external libraries for this question; it should be your own implementation.)
- (h) Using the training data for all digits, perform 3 fold cross-validation on your K-NN classifier and report your average accuracy.
- (i) Generate a confusion matrix (of size  $10 \times 10$ ) from your results. Which digits are particularly tricky to classify?
- (j) Train your classifier with all of the training data, and test your classifier with the test data. Submit your results to Kaggle.

### 2. The Titanic Disaster

- (a) Join the [Titanic: Machine Learning From Disaster](#) competition on Kaggle. Download the training and test data.
- (b) Using logistic regression, try to predict whether a passenger survived the disaster. You can choose the features (or combinations of features) you would like to use or ignore, provided you justify your reasoning.
- (c) Train your classifier using all of the training data, and test it using the testing data. Submit your results to Kaggle.

## WRITTEN EXERCISES

1. Variance of a sum. Show that the variance of a sum is  $\text{var}[X - Y] = \text{var}[X] + \text{var}[Y] - 2\text{cov}[X, Y]$ , where  $\text{cov}[X, Y]$  is the covariance between random variables  $X$  and  $Y$ .
2. Bayes rule for quality control. You're the foreman at a factory making ten million widgets per year. As a quality control step before shipment, you create a detector that tests for defective widgets before sending them to customers. The test is uniformly 95% accurate, meaning that the probability of testing positive given that the widget is defective is 0.95, as is the probability of testing negative given that the widget is not defective. Further, only one in 100,000 widgets is actually defective.
  - (a) Suppose the test shows that a widget is defective. What are the chances that it's actually defective given the test result?
  - (b) If we throw out all widgets that are defective, how many good widgets are thrown away per year? How many bad widgets are still shipped to customers each year?
3. In  $k$ -nearest neighbors, the classification is achieved by majority vote in the vicinity of data. Suppose our training data comprises  $n$  data points with two classes, each comprising exactly half of the training data, with some overlap between the two classes.
  - (a) Describe what happens to the average 0-1 prediction error on the training data when the neighbor count  $k$  varies from  $n$  to 1. (In this case, the prediction for training data point  $x_i$  includes  $(x_i, y_i)$  as part of the example training data used by  $k$ NN.)
  - (b) We randomly choose half of the data to be removed from the training data, train on the remaining half, and test on the held-out half. Predict and explain with a sketch how the average 0-1 prediction error on the held-out validation set might change when  $k$  varies? Explain your reasoning.
  - (c) We wish to choose  $k$  by cross-validation and are considering how many folds to use. Compare both the computational requirements and the validation accuracy of using different numbers of folds for  $k$ NN and recommend an appropriate number.
  - (d) In  $k$ NN, once  $k$  is determined, all of the  $k$ -nearest neighbors are weighted equally in deciding the class label. This may be inappropriate when  $k$  is large. Suggest a modification to the algorithm that avoids this caveat.
  - (e) Give two reasons why  $k$ NN may be undesirable when the input dimension is high.