
CS5785 Homework 1

This homework is due on **Friday, September 17th, 2021 at 11:59PM EST**. The homework is split into programming exercises and written exercises. Upload your homework to Gradescope. Your submission will have two parts:

1. A write-up as a single .pdf file. Submit this under the hw1-report assignment in Gradescope.
2. Source code and data files for all of your experiments (AND figures) in .py files if you use Python or .ipynb files if you use Jupyter Notebooks. 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.

These files should be made available via a **private** GitHub¹ repository (that you share with your teammate). Place all the relevant code files for this assignment in a hw1 folder in your repository and use the GitHub submission method under the hw1-code assignment in Gradescope. **Note:** Gradescope will take a snapshot of your GitHub repository during submission. Any changes to your repository after the assignment is submitted **will not** be reflected in your submission.

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 Canvas for announcements, policy changes, etc. and Slack for homework related questions. 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 get stuck, we encourage you to post a question on the *#homework* channel on Slack². That way, your solutions will be available to the other students in the class.
- The professor and TAs will offer office hours, 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).

²<https://join.slack.com/t/cs5785-2021fa-kallus/signup>

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 that implement K-NN for this question; it should be your own implementation. You can still use libraries such as numpy, scipy, pandas, etc., for data manipulation.)
Hint: You can use `scipy.spatial.distance.cdist`³ and `numpy.argsort`⁴ to make your implementation more efficient.
- (h) Randomly split the training data into two halves. Train your k-NN classifier on the first half of the data, and test it on the second half, reporting your average accuracy. *Note:* If you find your implementation to be slow, you can instead select a subsample of 5,000 from the test set rather than the entire test set.
- (i) Generate a confusion matrix (of size 10×10) 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. *Note:* You can also use scikit-learn's K-NN implementation⁵ for this.

2. The Titanic Disaster

- (a) Join the [Titanic: Machine Learning From Disaster](#) competition on Kaggle. Download the training and test data.

³<https://docs.scipy.org/doc/scipy/reference/generated/scipy.spatial.distance.cdist.html>

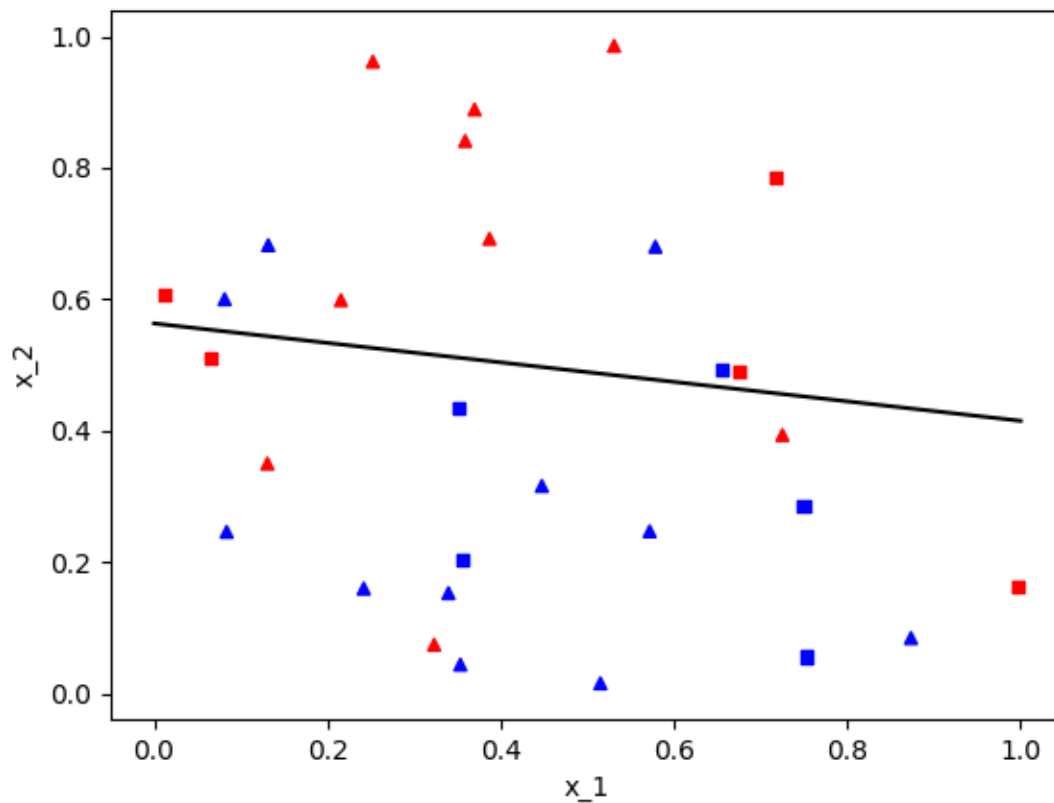
⁴<https://numpy.org/doc/stable/reference/generated/numpy.argsort.html>

⁵<https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html>

- (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 difference. Show that the variance of a difference 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 test as 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 plurality 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) 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.
 - (d) Give two reasons why k NN may be undesirable when the input dimension is high.
 - (e) The plot below shows a binary classification task with two classes: *red* and *blue*. Training data points are shown as triangles, and test data points are shown as squares, and all data points are colored according to their class. The line in the plot is the decision rule from training a logistic regression classifier on the training data, with points above the line classified as *red* and points below the line classified as *blue*.



- i. Calculate confusion matrices for both the training and test data for the logistic regression classifier.
- ii. Now, calculate confusion matrices for both the training and test data for the k NN algorithm with $k = 1$ on the same data, using the same train/test split (you should be able to do this by hand).
- iii. Discuss the difference in accuracy on both the training and test data between the two algorithms. Explain why this difference occurs.