

Due: Wednesday, January 26 at 11:59 pm

This homework is comprised of a set of coding exercises and a few math problems. While we have you train models across three datasets, the code for this entire assignment can be written in under 250 lines.

Start this homework early! You can submit to Kaggle only twice a day.

Deliverables:

1. Submit your predictions for the test sets to Kaggle as early as possible. Include your Kaggle scores in your write-up (see below).
2. Submit a **PDF** of your homework, **with an appendix listing all your code**, to the Gradescope assignment entitled “HW1 Write-Up”. You may typeset your homework in LaTeX or Word or submit neatly handwritten and scanned solutions. **Please start each question on a new page.** If there are graphs, include those graphs in the correct sections. **Do not** put them in an appendix. We need each solution to be self-contained on pages of its own.
 - On the first page of your write-up, please list students who helped you or whom you helped on the homework. (Note that sending each other code is not allowed.)
 - On the first page of your write-up, please copy the following statement and sign your signature next to it. (Mac Preview, PDF Expert, and FoxIt PDF Reader, among others, have tools to let you sign a PDF file.) We want to make *extra* clear the consequences of cheating.
“I certify that all solutions are entirely in my own words and that I have not looked at another student’s solutions. I have given credit to all external sources I consulted.”
3. Submit all the code needed to reproduce your results to the Gradescope assignment entitled “HW1 Code”. You must submit your code twice: once in your PDF write-up (above) so the readers can easily read it, and again in compilable/interpretable form so the readers can easily run it. Do **NOT** include any data files we provided. Please include a short file named README listing your name, student ID, and instructions on how to reproduce your results. Please take care that your code doesn’t take up inordinate amounts of time or memory.
The Kaggle score will not be accepted if the code provided a) does not compile or b) compiles but does not produce the file submitted to Kaggle.

Python Configuration and Data Loading

This section is only setup and has no requires no submitted solution. Please follow the instructions below to ensure your Python environment is configured properly, and you are able to successfully load the data provided with this homework. For all coding questions, we recommend using Anaconda for Python 3.

- (a) Either install Anaconda for Python 3, or ensure you're using Python 3. To ensure you're running Python 3, open a terminal in your operating system and execute the following command:

```
python --version
```

Do not proceed until you're running Python 3.

- (b) Install the following dependencies required for this homework by executing the following command in your operating system's terminal:

```
pip install scikit-learn scipy numpy matplotlib
```

Please use Python 3 with the modules specified above to complete this homework.

- (c) You will be running out-of-the-box implementations of Support Vector Machines to classify three datasets. You will find a set of .mat files in the data folder for this homework. Each .mat file will load as a Python dictionary. Each dictionary contains three fields:

- **training_data**, the training set features. Rows are sample points and columns are features.
- **training_labels**, the training set labels. Rows are sample points. There is one column: The labels corresponding to rows of training_data above.
- **test_data**, the test set features. Rows are sample points and columns are features. You will fit a model to predict the labels for this test set, and submit those predictions to Kaggle.

The three datasets for the coding portion of this assignment are described below.

- **mnist_data.mat** contains data from the MNIST dataset. There are 60,000 labeled digit images for training, and 10,000 digit images for testing. The images are flattened grayscale, 28×28 pixels. There are 10 possible labels for each image, namely, the digits 0–9.



Figure 1: Examples from the MNIST dataset.

- **spam_data.mat** contains featurized spam data. The labels are 1 for spam and 0 for ham. The data folder includes the script **featurize.py** and the folders **spam**, **ham** (not spam), and **test** (unlabeled test data); you may modify **featurize.py** to generate new features for the spam data.
- **cifar10_data.mat** contains data from the CIFAR10 dataset. There are 50,000 labeled object images for training, and 10,000 object images for testing. The images are flattened to $3 \times 32 \times 32$ (3 color channels). The labels 0–9 correspond alphabetically to the categories. For example, 0 means airplane, 1 means automobile, 2 means bird, and so on.

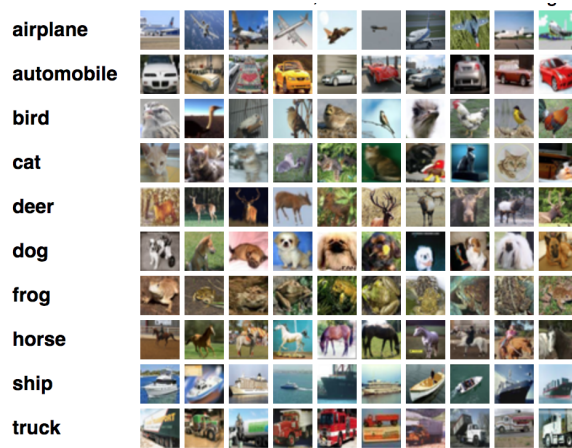


Figure 2: Examples from the CIFAR-10 dataset.

To check whether your Python environment is configured properly for this homework, ensure the following Python script executes without error. Pay attention to errors raised when attempting to import any dependencies. Resolve such errors by manually installing the required dependency (e.g. execute `pip install numpy` for import errors relating to the numpy package).

```
import sys
if sys.version_info[0] < 3:
    raise Exception("Python 3 not detected.")

import numpy as np
import matplotlib.pyplot as plt
from sklearn import svm
from scipy import io

for data_name in ["mnist", "spam", "cifar10"]:
    data = io.loadmat("data/%s_data.mat" % data_name)
    print("\nloaded %s data!" % data_name)
    fields = "test_data", "training_data", "training_labels"
    for field in fields:
        print(field, data[field].shape)
```

For the entire assignment, you may only use sklearn for the SVM model and the accuracy score function. Everything else (generating plots) must be done without the use of sklearn.

1 Data Partitioning

Rarely will you receive “training” data and “validation” data; usually you will have to partition available labeled data yourself. In this question, you will **shuffle and partition** each of the datasets in the assignment. Shuffling prior to splitting crucially ensures that all classes are represented in your partitions.

- (a) For the MNIST dataset, write code that sets aside 10,000 training images as a validation set.
- (b) For the spam dataset, write code that sets aside 20% of the training data as a validation set.
- (c) For the CIFAR-10 dataset, write code that sets aside 5,000 training images as a validation set.

Note: Make sure that you shuffle the labels with the training images. It’s a very common error to mislabel the training images by forgetting to permute the labels with the images!

Deliverable: Attach a copy of your data partitioning code to your homework report under question 1.

2 Support Vector Machines: Coding

We will use linear Support Vector Machines to classify our datasets. For images, we will use the simplest of features for classification: raw pixel brightness values. In other words, our feature vector for an image will be a row vector with all the pixel values concatenated in a row major (or column major) order.

There are several ways to evaluate models. We will use *classification accuracy*, or the percent of examples classified correctly, as a measure of the classifier performance (see here: https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html). Error rate, or one minus the accuracy, is another common metric.

Train a linear Support Vector Machine (SVM) on all three datasets. For each dataset, plot the accuracy on the training and validation sets versus the number of training examples that you used to train your classifier. The number of training examples to use are listed for each dataset in the following parts.

You may only use `sklearn` for the SVM model and the `sklearn.metrics.accuracy_score` function. Everything else (generating plots) must be done without the use of `sklearn`.

- (a) For the **MNIST** dataset, use raw pixels as features. Train your model with the following numbers of training examples: 100, 200, 500, 1,000, 2,000, 5,000, 10,000. At this stage, you should expect accuracies between 70% and 90%.

Hint: Be consistent with any preprocessing you do. Use either integer values between 0 and 255 or floating-point values between 0 and 1. Training on floats and then testing with integers is bound to cause trouble.

- (b) For the **spam** dataset, use the provided word frequencies as features. In other words, each document is represented by a vector, where the i^{th} entry denotes the number of times word i (as specified in `featurize.py`) is found in that document. Train your model with the following numbers of training examples: 100, 200, 500, 1,000, 2,000, **ALL**.

At this stage, you should expect accuracies between 70% and 90%.

- (c) For the **CIFAR-10** dataset, use raw pixels as features. At this stage, you should expect accuracies between 25% and 35%. Be forewarned that training SVMs for CIFAR-10 takes a couple minutes to run for a large training set locally. Train your model with the following numbers of training examples: 100, 200, 500, 1,000, 2,000, 5,000.

Note: You can use either `SVC(kernel='linear')` or `LinearSVC` as your SVM model, though they each solve slightly different optimization problems using different libraries. On MNIST, `LinearSVC` was faster on one member of Course Staff's laptop, though the exact results will likely depend on your computer, the parameters of the algorithm, and the data (number of data pts vs number of features).

Deliverable: For this question, you should include three plots showing number of examples versus training and validation accuracy for each of the datasets. Additionally, be sure to include your code in the "Code Appendix" portion of your write-up.

3 Hyperparameter Tuning

In the previous problem, you learned parameters for a model that classifies the data. Many classifiers also have *hyperparameters* that you can tune to influence the parameters. In this problem, we'll determine good values for the regularization parameter C in the soft-margin SVM algorithm.

When we are trying to choose a hyperparameter value, we train the model repeatedly with different hyperparameters. We select the hyperparameter that gives the model with the highest accuracy on the validation dataset. Before generating predictions for the test set, the model should be retrained using all the labeled data (including the validation data) and the previously-determined hyperparameter.

The use of automatic hyperparameter optimization libraries is prohibited for this part of the homework.

Deliverable: For the MNIST dataset, find the best C value. In your report, list at least 8 C values you tried, the corresponding accuracies, and the best C value. As in the previous problem, for performance reasons, you are required to train with up to 10,000 training examples but not required to train with more than that. Again, reference any code you used to perform a hyper-parameter sweep in the code appendix.

4 K-Fold Cross-Validation

For smaller datasets (e.g., the spam dataset), the validation set contains fewer examples, and our estimate of our accuracy might not be accurate—the estimate has high variance. A way to combat this is to use *k-fold cross-validation*.

In *k*-fold cross-validation, the training data is shuffled and partitioned into *k* disjoint sets. Then the model is trained on $k - 1$ sets and validated on the k^{th} set. This process is repeated *k* times with each set chosen as the validation set once. The cross-validation accuracy we report is the accuracy averaged over the *k* iterations.

Use of automatic cross-validation libraries is prohibited for this part of the homework.

Deliverable: For the spam dataset, use 5-fold cross-validation to find and report the best *C* value. In your report, list at least 8 *C* values you tried, the corresponding accuracies, and the best *C* value. Again, please include your code for cross validation or include a reference to its location in your code appendix.

Hint: Effective cross-validation requires choosing from **random** partitions. This is best implemented by randomly shuffling your training examples and labels, then partitioning them by their indices.

5 Kaggle

- MNIST Competition: <https://www.kaggle.com/c/hw1-mnist-competition-cs189sp22>
- SPAM Competition: <https://www.kaggle.com/c/hw1-spam-competition-cs189sp22>
- CIFAR-10 Competition: <https://www.kaggle.com/c/hw1-cifar-10-competition-cs189sp22>

With the best model you trained for each dataset, generate predictions for the test sets we provide and save those predictions to .csv files. Be sure to use integer labels (not floating-point!) and no spaces (not even after the commas). Upload your predictions to the Kaggle leaderboards (submission instructions are provided within each Kaggle competition). In your write-up, include your Kaggle name as it displays on the leaderboard and your Kaggle score for each of the three datasets.

General comments about the Kaggle sections of homeworks. Most or all of the coding homeworks will include a Kaggle section. Whereas other parts of a coding assignment might impose strict limits about what methods you're permitted to use, the Kaggle portions permit you to apply your creativity and find clever ways to improve your leaderboard performance. The main restriction is that you cannot use an entirely different learning technique. For example, this is an SVM homework, so you must use an SVM; you are not permitted to use a neural network or a decision tree instead of (or in addition to) a support vector machine. (You are also not allowed to search for the labeled test data and submit that to Kaggle; that's outright cheating.)

For example, to achieve higher positions on the Kaggle leaderboards, you may optionally add more features or use a nonlinear SVM kernel. Spam is a particularly good dataset for playing with feature engineering; one easy way to perform better in spam/ham is to add extra features with `featurize.py` (see below). Other examples of things you might investigate include SIFT and HOG features for images, and a bag-of-words model for spam/ham. For reasons we'll learn later this semester, dropping features that have little or no predictive power will often improve your test performance as much as adding the right new features. Although extensive creativity isn't generally necessary to get full points on an assignment, topping the Kaggle leaderboard gives your professor good material for letters of recommendation.

Whatever creative ideas you apply, please explain what you did in your write-up. Cite any external sources where you got ideas. If you have any questions about whether something is allowed or not, ask on Piazza.

Remember to start early! Kaggle only permits two submissions per leaderboard per day. To help you format the submission, please use `check.py` to run a basic sanity check on your submission and `save_csv.py` to help save your results.

To check your submission csv,

```
python check.py <competition name, eg. mnist> <submission csv file>
```

Deliverable: Your deliverable for this question has three parts. First, submit to all the Kaggle competitions listed above, and include your Kaggle score in your write-up. Second, include an explanation of what you tried, what worked, and what didn't to improve your accuracy. Finally, make sure to include all the code you used in the code appendix and provide a reference to it.

Modifying features for spam: The Python script `data/featurize.py` extracts features from the original emails in the Spam dataset. The spam emails can be found in `data/spam/`, the ham (ie. not spam) emails can be found in `data/ham/`, and the emails for the test set can be found in `data/test/`. You are encouraged to look at the emails and try to think of features you think would be useful in classifying an email as spam or ham.

To add a new feature, modify `featurize.py`. You are free to change the structure of the code provided, but if you are following the given structure, you need to do two things:

- Define a function, eg. `my_feature(text, freq)` that computes the value of your feature for a given email. The argument `text` contains the raw text of the email; `freq` is a dictionary containing the counts of each word in the email (or 0 if the word is not present). The value you return should be an integer or a float.
- Modify `generate_feature_vector` to append your feature to the feature vector. For example:

```
feature.append(my_feature(text, freq))
```

Once you are done modifying `featurize.py`, re-generate the training and test data by running

```
python data/featurize.py
```

6 Theory of Hard-Margin Support Vector Machines

A *decision rule* (or *classifier*) is a function $r : \mathbb{R}^d \rightarrow \pm 1$ that maps a feature vector (test point) to +1 (“in class”) or -1 (“not in class”). The decision rule for linear SVMs is of the form

$$r(x) = \begin{cases} +1 & \text{if } w \cdot x + \alpha \geq 0, \\ -1 & \text{otherwise,} \end{cases} \quad (1)$$

where $w \in \mathbb{R}^d$ and $\alpha \in \mathbb{R}$ are the parameters of the SVM. The primal hard-margin SVM optimization problem (which chooses the parameters) is

$$\min_{w, \alpha} |w|^2 \quad \text{subject to } y_i(X_i \cdot w + \alpha) \geq 1, \quad \forall i \in \{1, \dots, n\}, \quad (2)$$

where $|w| = \|w\|_2 = \sqrt{w \cdot w}$.

We can rewrite this optimization problem by using Lagrange multipliers to eliminate the constraints. (If you’re curious to know what Lagrange multipliers are, the Wikipedia page is recommended, but you don’t need to understand them to do this problem.) We thereby obtain the equivalent optimization problem

$$\max_{\lambda_i \geq 0} \min_{w, \alpha} |w|^2 - \sum_{i=1}^n \lambda_i (y_i (X_i \cdot w + \alpha) - 1). \quad (3)$$

Note: λ_i must be greater than or equal to 0.

(a) Show that Equation (3) can be rewritten as the *dual optimization problem*

$$\max_{\lambda_i \geq 0} \sum_{i=1}^n \lambda_i - \frac{1}{4} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j X_i \cdot X_j \quad \text{subject to } \sum_{i=1}^n \lambda_i y_i = 0. \quad (4)$$

Hint: Use calculus to determine and prove what values of w and α optimize Equation (3). Explain where the new constraint comes from.

SVM software usually solves this dual optimization program, not the primal optimization program.

(b) Suppose we know the values λ_i^* and α^* that optimize Equation (3). Show that the decision rule specified by Equation (1) can be written

$$r(x) = \begin{cases} +1 & \text{if } \alpha^* + \frac{1}{2} \sum_{i=1}^n \lambda_i^* y_i X_i \cdot x \geq 0, \\ -1 & \text{otherwise.} \end{cases} \quad (5)$$

(c) Applying Karush–Kuhn–Tucker (KKT) conditions (See Wikipedia for more information), any pair of optimal primal and dual solutions w^* , α^* , λ^* for a linear, hard-margin SVM must satisfy the following condition:

$$\lambda_i^* (y_i (X_i \cdot w^* + \alpha^*) - 1) = 0 \quad \forall i \in \{1, \dots, n\}$$

This condition is called *complementary slackness*. Explain what this implies for points corresponding to $\lambda_i^* > 0$. What relationship do they have with the margin?

- (d) The training points X_i for which $\lambda_i^* > 0$ are called the *support vectors*. In practice, we frequently encounter training data sets for which the support vectors are a small minority of the training points, especially when the number of training points is much larger than the number of features. Explain why the support vectors are the only training points needed to evaluate the decision rule.
- (e) Assume the training points X_i and labels y_i are linearly separable. Using the original SVM formulation (not the dual) prove that there is at least one support vector for each class, +1 and -1.

Hint: Use contradiction. Construct a new weight vector $w' = w/(1 + \epsilon/2)$ and corresponding bias α' where $\epsilon > 0$. It is up to you to determine what ϵ should be based on the contradiction. If you provide a symmetric argument, you need only provide a proof for one of the two classes.