

B551 Assignment 4: Machine learning

Fall 2018

Due: Sunday December 9, 11:59PM

(You may submit up to 48 hours late for a 10% penalty.)

This assignment will give you a chance to implement machine learning techniques on realistic classification problems.

You'll work in a group of 1-3 people for this assignment. You can choose the group yourself, and please fill out the form (sent out as an Announcement) so that we can create your repository. There is no starter code for this assignment so you can get started even while waiting for your repo to be set up. You should only submit **one** copy of the assignment for your team, through GitHub, as described below. All people on the team will receive the same grade on the assignment, except in unusual circumstances; we will collect feedback about how well your team functioned in order to detect these circumstances. Please read the instructions below carefully; we cannot accept any submissions that do not follow the instructions given here. Most importantly: please **start early**, and ask questions on Piazza or in office hours.

The following problems require you to write programs in Python. You may import standard Python modules for routines not related to AI, such as basic sorting algorithms and basic data structures like queues. You must write all of the rest of the code yourself. If you have any questions about this policy, please ask us. Make sure that your code works on the CS Linux machines before submission because that is where we will test it.

For each programming problem, please include a detailed comments section at the top of your code that describes: (1) a description of how you formulated the problem; (2) a brief description of how your program works; (3) a discussion of any problems, assumptions, simplifications, and/or design decisions you made; and (4) answers to any questions asked below in the assignment.

Academic integrity. You and your teammates may discuss the assignment with other people at a high level, e.g. discussing general strategies to solve the problem, talking about Python syntax and features, etc. You may also consult printed and/or online references, including books, tutorials, etc., but you must cite these materials (e.g. in source code comments). However, the work and code that you and your partners submit must be your group's own work, which your group personally designed and wrote. You may not share written answers or code with any other students except your own partner, nor may you possess code written by another student who is not your partner, either in whole or in part, regardless of format.

Image classification

In this assignment we'll study a straightforward image classification task. These days, all modern digital cameras include a sensor that detects which way the camera is being held when a photo is taken. This meta-data is then included in the image file, so that image organization programs know the correct orientation — i.e., which way is “up” in the image. But for photos scanned in from film or from older digital cameras, rotating images to be in the correct orientation must typically be done by hand.

Your task in this assignment is to create a classifier that decides the correct orientation of a given image, as shown in Figure 1.

Data. To help you get started, we've prepared a dataset of images from the Flickr photo sharing website. The images were taken and uploaded by real users from across the world, making this a challenging task on a very realistic dataset.¹ Since image processing is beyond the scope of this class, we don't expect you

¹All images in this dataset are licensed under Creative Commons, and there are no copyright issues with using them for classroom purposes under fair use guidelines. However, copyrights are still held by the photographers in most cases; you should not post these images online or use these images for any commercial purposes without checking with the original photographer, who can be identified using the Flickr URL described below.

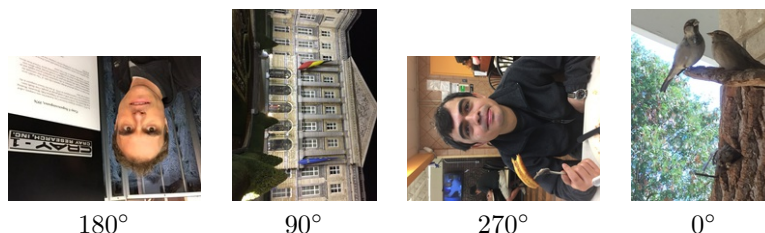


Figure 1: Some sample images, and their correct orientation labels.

to implement any special techniques related to images in particular. Instead, we'll simply treat the raw images as numerical feature vectors, on which we can then apply standard machine learning techniques. In particular, we'll take an $n \times m \times 3$ color image (the third dimension is because color images are stored as three separate planes – red, green, and blue), and append all of the rows together to produce a single vector of size $1 \times 3mn$.

We've done this work for you already, so that you can treat the images simply as vectors and do not have to worry about them being images at all. Thorough Canvas we've provided two ASCII text files, one for the training dataset and one for testing, that contain the feature vectors. To generate this file, we rescaled each image to a very tiny “micro-thumbnail” of 8×8 pixels, resulting in an $8 \times 8 \times 3 = 192$ dimensional feature vector. The text files have one row per image, where each row is formatted like:

```
photo_id correct_orientation r11 g11 b11 r12 g12 b12 ...
```

where:

- **photo_id** is a photo ID for the image.
- **correct_orientation** is 0, 90, 180, or 270. Note that some small percentage of these labels may be wrong because of noise; this is just a fact of life when dealing with data from real-world sources.
- **r11** refers to the red pixel value at row 1 column 1, **r12** refers to red pixel at row 1 column 2, etc., each in the range 0-255.

Although you can get away with just using the above text files, you may want to inspect the original images themselves, for debugging or analysis purposes, or if you want to change something about the way we've created the feature vectors (e.g. experiment with larger or smaller “micro-thumbnails”). You can view the images in two different ways:

- You can view the original high-resolution image on Flickr.com by taking just the numeric portion of the photo_id in the file above (e.g. if the photo_id in the file is `test/123456.jpg`, just use 123456), and then visiting the following URL:
http://www.flickr.com/photo_zoom.gne?id=numeric_photo_id
- We'll provide a zip file of all the images in JPEG format on Canvas. We've reduced the size of each image to a 75×75 square, which still (usually) preserves enough information to figure out the image orientation. The ZIP file also includes UNIX scripts that will convert images in the zip file to the ASCII feature vector file format above. If you want, this lets you experiment with modifying the script to produce other feature vectors (e.g. smaller sized images, or in different color mappings, etc.) and to run your classifier on your own images.

The training dataset consists of about 10,000 images, while the test set contains about 1,000. For the training set, we've rotated each image 4 times to give four times as much training data, so there are about 40,000

lines in the train.txt file (the training images on Flickr and the ZIP file are all oriented correctly already). In test.txt, each image occurs just once (rotated to a random orientation) so there are only about 1,000 lines. If you view these images on Flickr, they'll be oriented correctly, but those in the ZIP file may not be.

What to do. Your goal is to implement and test several different classifiers on this problem: k -nearest neighbors, AdaBoost, and decision forests. For training, your program should be run like this:

```
./orient.py train train_file.txt model_file.txt [model]
```

where [model] is one of **nearest**, **adaboost**, **forest**, or **best**. This program uses the data in train_file.txt to produce a trained classifier of the specified type, and save the parameters in model_file.txt. You may use any file format you'd like for model_file.txt; the important thing is that your test code knows how to interpret it.

For testing, your program should be run like this:

```
./orient.py test test_file.txt model_file.txt [model]
```

where [model] is again one of **nearest**, **adaboost**, **forest**, **best**. This program should load in the trained parameters from model_file.txt, run each test example through the model, display the classification accuracy (in terms of percentage of correctly-classified images), and output a file called output.txt which indicates the estimated label for each image in the test file. The output file should correspond to one test image per line, with the photo_id, a space, and then the estimated label, e.g.:

```
test/124567.jpg 180
test/8234732.jpg 0
```

Here are more detailed specifications for each classifier.

- **nearest:** At test time, for each image to be classified, the program should find the k “nearest” images in the training file, i.e. the ones with the closest distance (least vector difference) in Euclidean space, and have them vote on the correct orientation.
- **adaboost:** Use very simple decision stumps that simply compare one entry in the image matrix to another, e.g. compare the red pixel at position 1,1 to the green pixel value at position 3,8. You can try all possible combinations (roughly 192^2) or randomly generate some pairs to try.
- **forest:** A decision forest is a generalization of a decision tree: instead of constructing a single decision tree for a given classification problem, the idea is to instead create many – dozens, hundreds, or even thousands – of decision trees, and then take a vote of these trees to make a decision. Each of the trees is learned with a *different random subset of the training data*, so that each tree becomes an “expert” on different subsets of images, and the forest together forms a “committee of experts.” This is a purposely open-ended problem because you have to decide how to define the features, as well as how deep to make the trees and how many trees to learn. (As a starting point, you might use features that simply compare individual pixel values with a threshold – e.g., whether the red pixel at position 5,5 is greater than or less than a threshold τ .)
- **best:** Use whichever algorithm (one of the above, or something completely different!) that you recommend to give the best accuracy.

Note that, you have to implement these classifiers from scratch. It's not allowed to use any pre-implemented packages, e.g., scikit-learn. Please check with us before using packages other than os, sys, math, matplotlib, scipy, and numpy.

Each of the above machine learning techniques has a number of parameters and design decisions. It would be impossible to try all possible combinations of all of these parameters, so identify a few parameters and

conduct experiments to study their effect on the final classification accuracy. In your report, present neatly-organized tables or graphs showing classification accuracies and running times as a function of the parameters you choose. Which classifiers and which parameters would you recommend to a potential client? How does performance vary depending on the training dataset size, i.e. if you use just a fraction of the training data? Show a few sample images that were classified correctly and incorrectly. Do you see any patterns to the errors?

As in the past, a small percentage of your assignment grade will be based on how accurate your “best” algorithm is with respect to the rest of the class. We will use a separate test dataset, so make sure to avoid overfitting!

Although we expect the testing phase of your classifiers to be relatively fast, we will not evaluate the efficiency of your training program. Please commit your model files for each method to your git repo when you submit your source code, and please name them `nearest_model.txt`, `adaboost_model.txt`, `forest_model.txt`, and `best_model.txt`.

What to turn in

Turn in your source code file(s) and your report by pushing to GitHub (remember to **add**, **commit**, **push**) — we’ll grade whatever version you’ve put there as of 11:59PM on the due date. To make sure that the latest version of your work has been accepted by GitHub, you can log into the github.iu.edu website and browse the code online. Make sure to include your trained model files!