

CSci 4270 and 6270
Computational Vision,
Spring Semester 2021
Homework 5
Due: Friday, April 9, at 11:59:59 pm EDT

Overview

The focus of this assignment is scene classification. In particular you will write SVM and neural network classifiers to determine the “background” class of a scene.

The five possibilities we will consider are grass, wheat field, road, ocean and red carpet. Some of these are relatively easy, but others are hard. A zip file containing the images can be found at

https://drive.google.com/file/d/1nrALsQXsU6KoPAMCz65H0S20N_mRvixZ/view?usp=sharing

The images are divided into three sets: training, validation and test. Training images are used directly to optimize the learned weights of a model. The validation set is used to tune the parameters that determine the model that is trained. These parameters are things like (a) the number of layers of a neural network, (b) its learning rate, (c) when to stop training, and (d) the parameter controlling the tradeoff between the margin and the slack values for a linear SVM. Once the parameters of the model are set and the model is trained, the model is applied to the test data just once to see how well it ultimately works.

Here’s a little more about the use of the validation data. Think of it as part of an outer loop of training where (a) the model parameters are set, (b) the model is trained, (c) the trained model is applied to the validation data to determine model accuracy, (d) if the resulting model is more accurate than any previously trained model, the model parameters and trained weights are saved, (e) new model parameters are set, and the process is repeated.

Classifiers to Build

In Problem 1 you will implement a descriptor and a linear SVM to classify the scene, while in Problem 2 you will implement two neural networks.

1. **(60 points)** In Lecture 14, which covered detection and SVMs, we focused on the “HoG” — histogram of oriented gradients — descriptor. After this method was published, many different types of descriptors were invented for many applications. For this problem you are going to implement a descriptor that combines location and color. It will include no gradient information. One long descriptor vector will be computed for each image, and then a series of SVM classifiers will be applied to the descriptor vector to make a decision. We strongly urge you to write two scripts to solve this problem, one to compute and save the descriptor for each image, and one to reload the descriptors, train the classifiers, and evaluate the performance.

The descriptor is formed from an image by computing a 3d color histogram in each of a series of overlapping subimages, unraveling each histogram into a vector (1d NumPy array), and concatenating the resulting vectors into one long vector. The key parameter here will be the value of t , the number of histogram bins in each of the three color dimensions. Fortunately, calculation of color histograms is straightforward. Here is example code to form an image of random values and compute its 3d color histogram:

```
import numpy as np
```

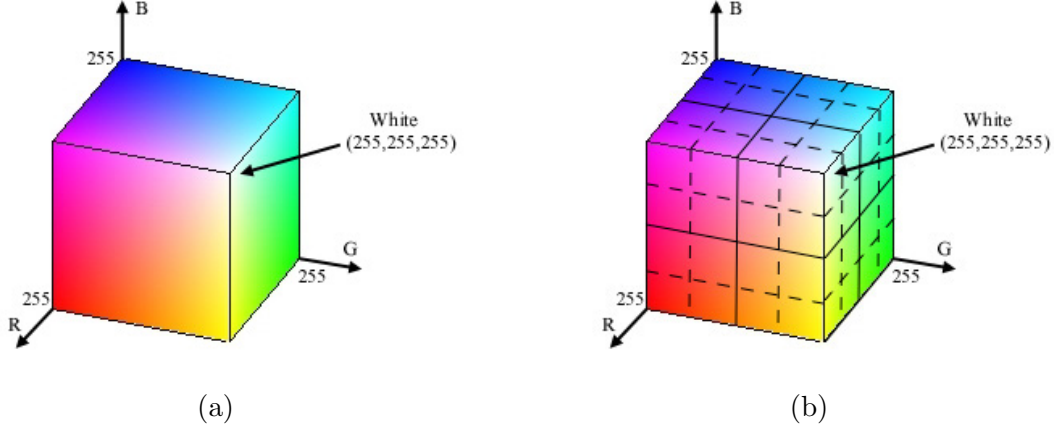


Figure 1: Color histogram bins

```
# Generate a small image of random values.
M, N = 20, 24
img = np.floor(256 * np.random.random(M * N * 3)).astype(np.uint8)

# Calculate and print the 3d histogram.
t = 4
pixels = img.reshape(M*N, 3)
hist, _ = np.histogramdd(pixels, (t, t, t))
print('histogram shape:', hist.shape)    # should be t, t, t
print('histogram:', hist)
print(np.sum(hist))                      # should sum to M*N
```

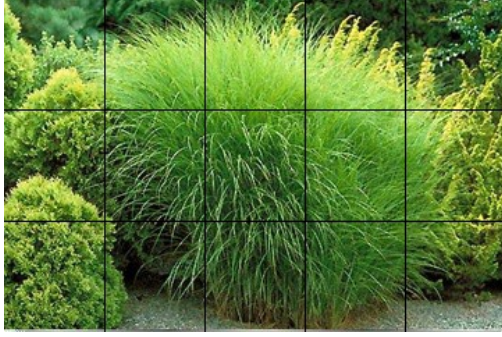
Each histogram bin covers a range of $u = 256/t$ gray levels in each color dimension. For example, $\text{hist}[i, j, k]$ is the number of pixels in the random image whose red value is in the range $[i * u, (i + 1) * u)$ **and** whose green value is in the range $[j * u, (j + 1) * u)$ **and** whose blue value is in the range $[k * u, (k + 1) * u)$. Figure 1 illustrates the partitioning of the RGB color cube into $4 \times 4 \times 4$ regions.

As mentioned above, the image will be divided into overlapping subimage blocks and one color histogram size of t^3 will be computed in each block. These will be concatenated to form the final histogram. Let b_w be the number of subimage blocks across the image and let b_h be the number of blocks going down the image. This will produce $b_w \cdot b_h$ blocks overall, and a final descriptor vector of size $t^3 \cdot b_w \cdot b_h$. To compute the blocks in an image of $H \times W$ pixels, let

$$\Delta_w = \frac{W}{b_w + 1} \quad \text{and} \quad \Delta_h = \frac{H}{b_h + 1}.$$

The blocks will each cover $2\Delta_w$ columns and $2\Delta_h$ rows of pixels. The image pixel positions of the upper left corners of the blocks will be

$$(m\Delta_w, n\Delta_h), \quad \text{for } m = 0, \dots, b_w - 1, \quad n = 0, \dots, b_h - 1.$$



Original image with Δ_w and Δ_h spaced lines Blocks of pixels over which histograms are formed

Figure 2: Image block tiling for $b_w = 4$ and $b_h = 2$.

Note that some pixels will contribute to only one histogram, some will contribute to two, and others will contribute four. (The same is true of the HoG descriptor.) Figure 2 illustrates the formation of blocks.

After you have computed the descriptors, you will train a series of SVM classifiers, one for each class. To do so, you will be given a set of 4000 training images, $\{I_i\}$, with class labels $y_i \in (1, \dots, k)$ (for us, $k = 5$). To train classifier C_j , images with label $y_i = j$ are treated as $y_i = +1$ in linear SVM training and images with label $y_i \neq j$ are treated as $y_i = -1$. This will be repeated for each of the k classifiers. The descriptor is computed for each training image I_i to form the data vectors \mathbf{x}_i .

Each resulting classifier C_j will have a weight vector \mathbf{w}_j and offset b_j . The score for classifier j for a test image with descriptor vector \mathbf{x} is

$$d_j = \frac{1}{\|\mathbf{w}_j\|} [\mathbf{w}_j^\top \mathbf{x} + b_j].$$

(Recall that the $1/\|\mathbf{w}_j\|$ ensures that d_j is a signed distance.) The classification for the test image I is the class associated with the value of j that gives the maximum d_j score. This is used even if none of the d_j scores are positive.

For each SVM classifier, output simple statistics on the validation steps you use to set the model parameters. This should likely just be the multiplier that controls tradeoff between margin and slack values. Note that a different value of the multiplier is allowed for each of the five SVMs you train. Reasonable ranges for this parameter are $[0.1, 10]$

After you complete training, you will test your classifiers with the set of 750 test images. Each will be run as described above and the label will be compared to the known correct label. You will output the percentage correct for each category, followed by a $k \times k$ confusion matrix. The confusion matrix entry at row r and column c shows the number of times when r was the correct class label and c was the chosen class label. The confusion matrix would have all 0's in the non-diagonal entries when the SVM classifier is operating at 100% accuracy.

Some Details

- (a) It is very important that you use `histogramdd` or a similar function to compute the histogram. Do not iterate over the pixels in the image using for loops: you will have serious trouble handling the volume of images. You can iterate over the indices of the subimage blocks and then use Numpy calls as above to form the histograms.
- (b) As mentioned above, we suggest that you write one script to compute and save the descriptors for all training and test images, and then write a second script to train your SVM classifiers and generate your test output. We suggest that you use Python's `pickle` module.
- (c) For your SVM implementation, we suggest using `sklearn.svm.LinearSVC`. To use the scikit-learn (sklearn) Python module, you will need to install the package. If you are using Anaconda, as suggested earlier in the semester, you can simply run

`conda install scikit-learn`
- (d) The confusion matrix can be made using Matplotlib or
`sklearn.metrics.confusion_matrix`
- (e) The computation of feature extraction might still be time consuming even with efficient Numpy use. We suggest that you develop and debug your program using a subset of the training image set before running your final version on the full training and test sets.
- (f) We suggest using at a minimum $t = 4$ and $b_w = b_h = 4$. This will give a descriptor of size 1,024 per image.
- (g) Finally, you can investigate the addition of gradient histograms to your descriptors. Doing so, with a careful analysis of the impacts, can earn you up to 5 points extra credit.

Submit your code, your output showing your validation and training experiments, your final test results, and a write-up describing design choices, your validation steps and a summary discussion of when your classifier works well, when it works poorly, and why.

2. **(60 points)** We continue the background classification problem using neural networks. Specifically, you will use *pytorch* to implement two neural networks, one using only fully-connected layers, and the other using convolutional layers in addition to fully-connected layers. The networks will each start directly with the input images so you will not need to write or use any code to do manual preprocessing or descriptor formation. Therefore, once you understand how to use *pytorch* the code you actually write will in fact be quite short. Your output should be similar in style to your output from Problem 1.

Make sure that your write-up includes a discussion of the design of your networks and your validation steps.

To help you get started I've provided a Jupyter notebook (on the Submittity site) that illustrates some of the main concepts of *pytorch*, starting with Tensors and Variables and proceeding to networks, loss functions, and optimizers. This also includes pointers to tutorials on pytorch.org. This notebook will be discussed in class on March 29 and April 1. If you already know TensorFlow, you'll find *pytorch* quite straightforward.

Two side notes:

- (a) PyTorch includes some excellent tools for uploading and transforming images into Tensor objects. For this assignment, you will not need to use these since you've already written

code for image input for the previous problem that gathers images into numpy arrays and it is trivial to transform these objects into pytorch tensors.

- (b) Because of the size of the images, you might find it quite computationally expensive and tedious to use a fully-connected network on the full-sized images. Therefore, you are welcome to resize the images before input to your first network. In fact, we strongly suggest that you start with significantly downsized images first and then see how far you can scale up!

Access to a GPU and to AIMOS

For your experiments it will be helpful to have access to a GPU. Many of your computers have them, but there are other possibilities. The simplest one is to use Google Colaboratory, and the Jupyter notebook distributed for lecture on Monday March 29 demonstrates how to do this. If you'd like more power or if you'd just like to explore, we have the chance to access the RPI AIMOS super computer <https://cci.rpi.edu/aimos>. Instructions for signing up for an account and for setting started will be posted on-line.