

HOMEWORK 1: SPATIAL PYRAMID MATCHING FOR SCENE CLASSIFICATION

16-720A Computer Vision (Spring 2023)

<https://canvas.cmu.edu/courses/32966>

OUT: Jan 25th, 2023

DUE: Feb 13th, 2023 11:59 PM

Instructor: Deva Ramanan

TAs: Kangle Deng, Vidhi Jain, Xiaofeng Guo, Chung Hee Kim, Ingrid Navarro Anaya



Figure 0.1: **Scene Classification:** Given an image, can a computer program determine where it was taken? In this homework, you will build a representation based on bags of visual words and use spatial pyramid matching for classifying the scene categories.

START HERE: Instructions

- Please refer to the [course logistics page](#) for information on the **Collaboration Policy** and **Late Submission Policy**.
- **Submitting your work:** There will be two submission slots for this homework on **Gradescope**: Written and Programming.
 - For written problems such as short answer, multiple choice, derivations, proofs, or plots, we will be using the written submission slot. Please use this provided template. **We don't accept handwritten submissions.** Each answer should be completed in the boxes provided below the question. You are allowed to adjust the size of these boxes, but **make sure to link your answer to each question when submitting to Gradescope.** Otherwise, your submission will not be graded.
 - You are also required to upload your code, which you wrote to solve this homework, to the Programming submission slot. Your code may be run by TAs so please make sure it is in a

workable state. The assignment must be completed using Python 3.7 or newer. We recommend setting up a [conda environment](#), but you are free to set up your environment however you like.

- Regrade requests can be made after the homework grades are released, however this gives the TA the opportunity to regrade your entire paper, meaning if additional mistakes are found then points will be deducted.
- **Start early!** This homework may take a long time to complete.
- **Attempt to verify your implementation as you proceed.** If you don't verify that your implementation is correct on toy examples, you will risk having a huge mess when you put everything together. Here are two tips:
 - (1) Once you write a function, uncomment the corresponding lines in `main.py` to verify whether the function executes correctly.
 - (2) To debug your logic within a function, use `print()` or `breakpoint()`.
- Follow the guidelines in Section 5: HW Checklist for writeup and code. If you have any questions or need clarifications, please post in Slack or visit the TAs during office hours.

Overview

Bag-of-words (BoW) can be applied to many problems in computer vision, including object recognition [5, 7] and scene classification [6, 8]¹. This homework will explore classic BoW along with extensions, such as pyramid matching [2, 4] and feature encoding [1]. Fig 0.2 provides an overview. Section 1 builds a dictionary of visual words from a training set of images by clustering. Section 2 builds a representation for a particular image as a histogram over visual words, or BoW. Finally, you will build a scene recognition system that classifies a test image by comparing it to a training library of images in BoW space (e.g., nearest-neighbor classification).

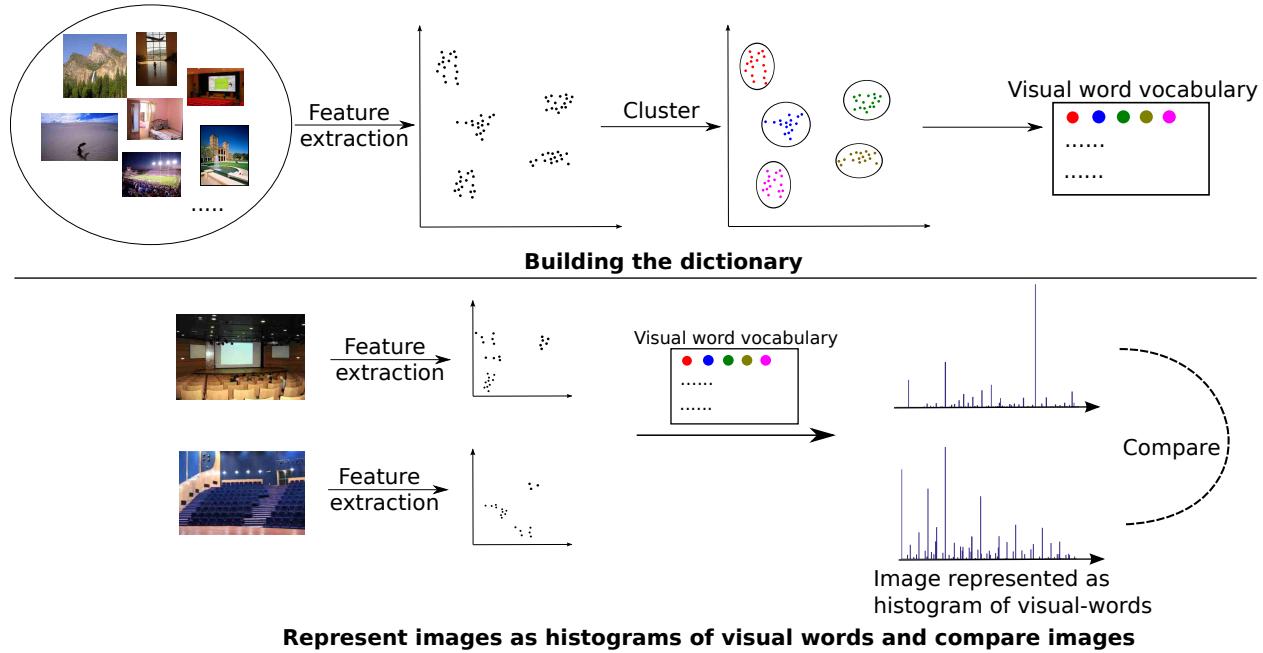


Figure 0.2: An overview of the bags-of-words approach to be implemented in the homework. First, given the training set of images, we extract the visual features of the images. In our case, we will use the filter responses of the pre-defined filter bank as the visual features. Next, we build visual words, *i.e.* a dictionary, by finding the centers of clusters of the visual features. To classify new images, we first represent each image as a vector of visual words, and then compare new images to old ones in the visual-word vector space – the nearest match provides a label!

What you will be doing: You will implement a scene classification system that uses the bag-of-words approach with its spatial pyramid extension. The paper that introduced the pyramid matching kernel [2] is

K. Grauman and T. Darrell. *The Pyramid Match Kernel: Discriminative Classification with Sets of Image Features*. ICCV 2005. http://www.cs.utexas.edu/~grauman/papers/grauman_darrell_iccv2005.pdf

Spatial pyramid matching [4] is presented in

S. Lazebnik, C. Schmid, and J. Ponce, *Beyond Bags of Features: Spatial Pyramid Matching for Recognizing Natural Scene Categories*, CVPR 2006. <http://www.di.ens.fr/willow/pdfs/cvpr06b.pdf>

¹This homework is largely self-contained, but reading the listed papers (or even just skimming them) will likely be helpful.

You will be working with a subset of the SUN database². The data set contains 1600 images from various scene categories like “aquarium”, “desert” and “kitchen”. And to build a recognition system, you will:

- take responses of a filter bank on images and build a dictionary of visual words, and then
- learn a model for images based on the bag of words (with spatial pyramid matching [4]), and use nearest-neighbor to predict scene classes in a test set.

In terms of number of lines of code, this assignment is fairly small. However, it may take *a few hours* to finish running the baseline system, so make sure you start early so that you have time to debug things. Try printing statements within long-running functions to verify that the function did not hang. Also, try **each component on a subset of the data set** first before putting everything together. We provide you with a number of functions and scripts in the hopes of alleviating some tedious or error-prone sections of the implementation. You can find a list of files provided in Section 4. *Though not necessary, you are recommended to implement a multi-processing³ version to make use of multiple CPU cores to speed up the code.* Functions with `n_workers` as input can benefit greatly from parallel processing.

Hyperparameters: We provide you with a basic set of hyperparameters, which might not be optimal. You will be asked in Q3.1 to tune the system you built and we suggest you to keep the defaults before you get to Q3.1. All hyperparameters can be found in a single configuration file `opts.py`.

²<http://groups.csail.mit.edu/vision/SUN/>

³Note that multi-threading in python does not make use of multiple CPU cores. It may not work on windows jupyter notebook.

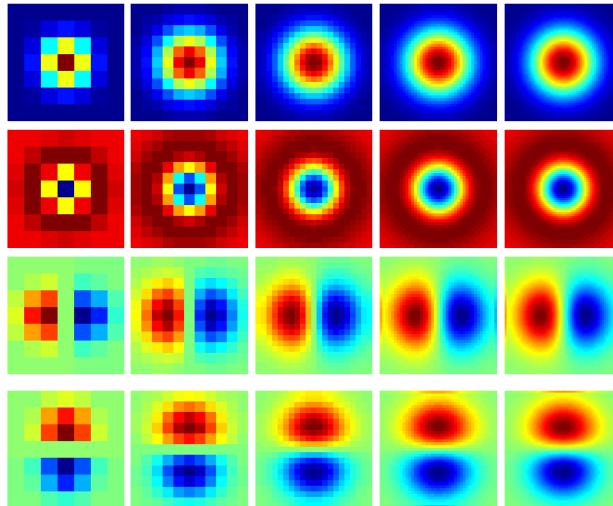


Figure 1.1: Multi-scale filter bank

1 Representing the World with Visual Words

1.1 Extracting Filter Responses

We want to run a filter bank on an image by convolving each filter in the bank with the image and concatenating all the responses into a vector for each pixel. In our case, we will be using 4 types of filters of multiple scales (`opts.filter_scales`). The filters are: (1) Gaussian, (2) Laplacian of Gaussian, (3) derivative of Gaussian in the x direction, and (4) derivative of Gaussian in the y direction.

Q1.1.1 (5 points): (a) What properties do each of the filter functions pick up? (See Fig 1.1) Try to group the filters into broad categories (e.g. all the Gaussians).

(b) Why do we need multiple scales of filter responses?

Q1.1.1 (a)(b)

Gaussian Filter: The Gaussian filter is a fuzzy filter that smooths/blurs out the input image. Unlike the median filter (non-linear filter for smoothing), the Gaussian filter preserves edge information while removing noise.

Laplacian of Gaussian Filter: It is a second-order derivative filter. When this filter is applied to the input image, we end up with zero crossings i.e. positions in the image where there is a sharp change in intensity (+ve to -ve or -ve to +ve through zero). In this case, however, we can also use this filter to detect blob-like features. A thing to note here is that the Laplace filter is noisy and hence to accurately detect the blobs, a gaussian filter is applied as a preprocessing step to smoothen the image.

Gaussian Filter (x derivative): First order derivative filter in the x direction which is used to calculate the vertical edges. Gaussian smoothing is applied to reduce noise

Gaussian Filter (y derivative): First order derivative filter in the y direction which is used to calculate the horizontal edges. Gaussian smoothing is applied to reduce noise

We use different scales primarily to extract different types of edges. For a larger scale, large-scale edges are detected and for smaller values, finer features are detected. We would like to include both types to build a robust classification system.

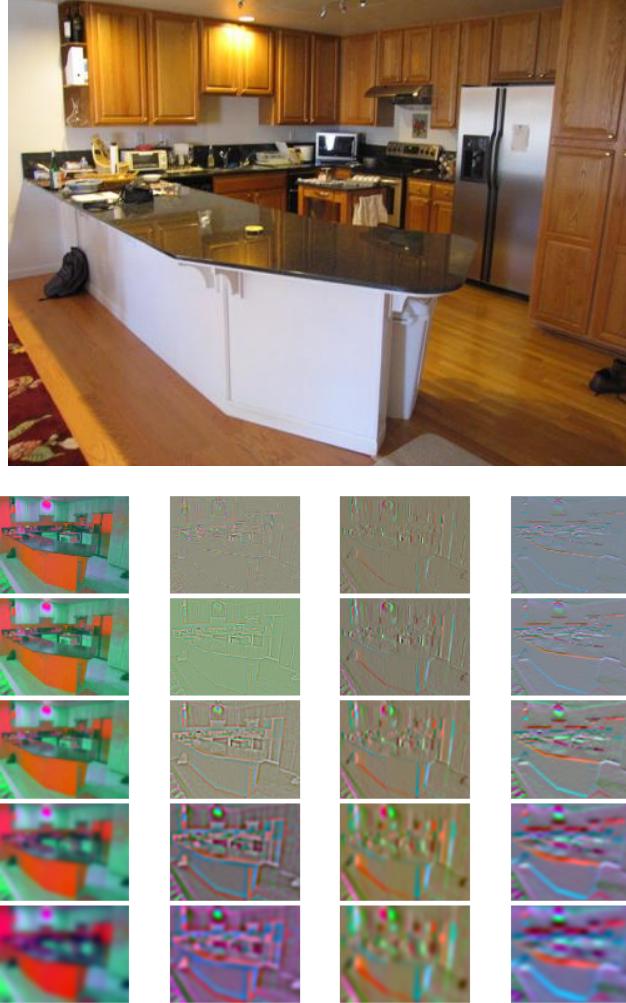


Figure 1.2: An input image and filter responses for all of the filters in the filter bank. **Top:** The input image. **Bottom:** The filter responses in Lab colorization, corresponding to the filters in Fig 1.1 (Transposed).

Q1.1.2 (10 points): For the code, loop through the filters and the scales to extract responses. Since color images have 3 channels, you are going to have a total of $3F$ filter responses per pixel if the filter bank is of size F . Note that in the given dataset, there are some gray-scale images. For those gray-scale images, you can simply duplicate them into three channels. Then output the result as a $3F$ channel image. Image laundromat/sun_afrrjykuhhlwiwun.jpg has 4 channels instead of 3. Discard the last channel. Try to first iterate across scales and then for each scale, iterate across each channel (i.e. $\text{Scale}_1 \{\text{Gaussian}\{\text{R,G,B}\}, \text{Laplacian}\{\text{R, G, B}\}, \dots\}$, $\text{Scale}_2 \{\text{Gaussian}\{\text{R,G,B}\}, \text{Laplace}\{\text{R, G, B}\}, \dots\}$). Use zero-padding if necessary. Normalize the input before passing the image to extract_filter_responses. Complete the function

```
visual_words.extract_filter_responses(opts, img)
```

and return the responses as `filter_responses`. We have provided you with template code, with detailed instructions commented inside. The convolution routine function `scipy.ndimage.convolve()` can be used with user-defined filters, but the functions `scipy.ndimage.gaussian_filter()` and `scipy.ndimage.gaussian_laplace()` may be useful here for improved efficiency. Note that by default `scipy.ndimage` applies filters to all dimensions including channels. Therefore you might want

to filter each channel separately. You can also pass in a parameter indicating you want either the x or y derivative.

Remember to check the input argument `image` to make sure it is a floating point type with range $[0, 1]$, and convert it if necessary. Be sure to check the number of input image channels and convert it to 3-channel if it is not. Before applying the filters, use the function `skimage.color.rgb2lab()` to convert your image into the Lab color space, which is designed to more effectively quantify color differences with respect to human perception. (See [here](#) for more information.) If the input `image` is an $M \times N \times 3$ matrix, then `filter_responses` should be a matrix of size $M \times N \times 3F$. Make sure your convolution function call handles image padding along the edges sensibly.

Apply all 4 filters at least 3 scales on `aquarium/sun_aztvjgubyrgvirup.jpg`, and visualize the responses as an image collage as shown in Fig 1.2. To plot the collage, you can use the included helper function `util.display_filter_responses` by providing a list of filter responses with those of the Lab channels grouped together with shape $M \times N \times 3$. We provide the skeleton code from line 17-21 in `main.py`. You can get the results by running `python main.py --filter-scales 1 2 4`.

Q1.1.2 Solution Collage of images

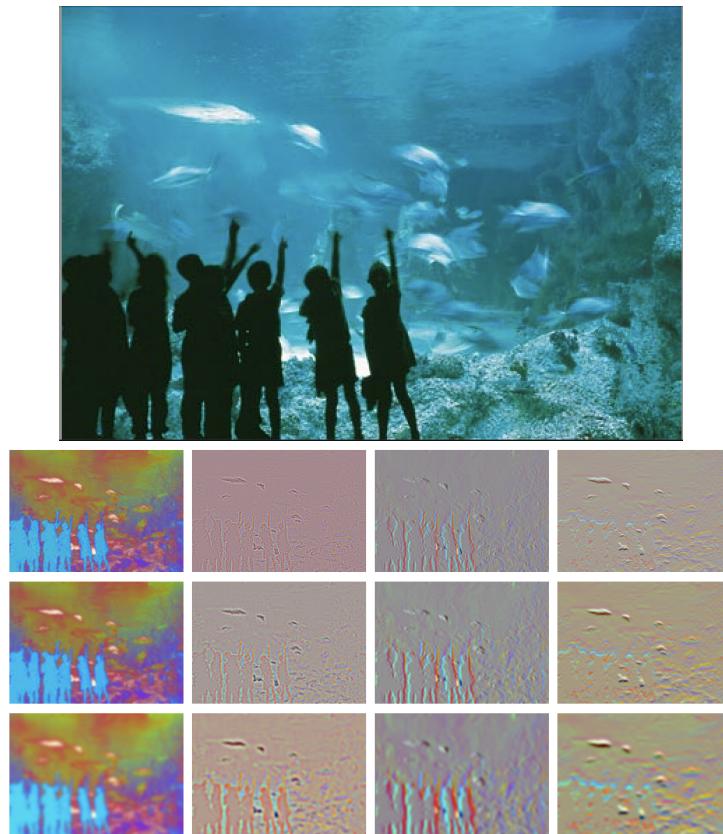


Figure 1.3: Input image (`aquarium/sun_aztvjgubyrgvirup.jpg`) and filter responses for all of the filters in the filter bank. **Top:** The input image. **Bottom:** The filter responses in Lab colorization (filter scales [1, 2, 4] where filter scale 1 has been applied to the first response and so on)

1.2 Creating Visual Words

You will now create a dictionary of visual words from the filter responses using k-means. After applying k-means, similar filter responses will be represented by the same visual word. You will use a dictionary with a fixed size. Instead of using all of the filter responses (**which might exceed the memory capacity of your computer**), you will use responses at α random pixels. If there are T training images, then you should collect a matrix `filter_responses` over all the images that is $\alpha T \times 3F$, where F is the filter bank size. Then, to generate a visual words dictionary with K words (`opts.K`), you will cluster the responses with k-means using the function `sklearn.cluster.KMeans` as follows:

```
kmeans = sklearn.cluster.KMeans(n_clusters=K).fit(filter_responses)
dictionary = kmeans.cluster_centers_
```

If you like, you can pass the `n_jobs` argument into the `KMeans()` object to utilize parallel computation.

Q1.2 (10 points): Write the functions

```
visual_words.compute_dictionary(opts, n_worker),
visual_words.compute_dictionary_one_image(args) (optional, multi-processing),
```

Given a dataset, these functions generate a dictionary. The overall goal of `compute_dictionary()` is to load the training data, iterate through the paths to the image files to read the images, and extract αT filter responses over the training files, and call k-means. This can be slow to run; however, the images can be processed independently and in parallel. Inside `compute_dictionary_one_image()`, you should read an image, extract the responses, and save to a temporary file. Here `args` is a collection of arguments passed into the function. Inside `compute_dictionary()`, you should load all the training data and create subprocesses to call `compute_dictionary_one_image()`. After all the subprocesses finish, load the temporary files back, collect the filter responses, and run k-means. A list of training images can be found in `data/train_files.txt`.

Finally, execute `compute_dictionary()`, and go do some push-ups while you wait for it to complete. If all goes well, you will have a file named `dictionary.npy` (with size of $K \times 3F$) that contains the dictionary of visual words. If the clustering takes too long, reduce the number of clusters and samples. You can start with a tiny subset of training images for debugging. We provide the skeleton code from line 24-25 in `main.py`. You can get the results by running `python main.py --filter-scales 1 2 4 --feat-dir TMP_OUT_DIR_FOR_EACH_IMG --out-dir FINAL_OUT_DIR`.

Include your implemented functions within the `minted` block below `compute_dictionary`, and optionally, `compute_dictionary_one_image` or other customized functions).

Q1.2

```
# Copy and paste your code here.

def compute_dictionary_one_image(args):
    """
    Extracts a random subset of filter responses of an image and save it to disk
    This is a worker function called by compute_dictionary

    You are free to make your own interface based on how you implement
    ← compute_dictionary
    """

    opts, data_dir_img, feat_dir_img = args
    with Image.open(data_dir_img) as img:
        img = np.array(img).astype(np.float32) / 255
        filter_responses = extract_filter_responses(opts, img)
        h, w, nf = filter_responses.shape
        # want 1 x alpha x 3F, 1 since only 1 image processed
        reshaped_filter_responses = np.reshape(filter_responses, (h * w, nf))
        rand_indices = random.sample(range(0, h * w), opts.alpha)
        filter_responses = reshaped_filter_responses[rand_indices, :]
        with open(feat_dir_img, "wb") as f:
            np.save(f, filter_responses)
    return feat_dir_img

def compute_dictionary(opts, n_worker=1):
    """
    Creates the dictionary of visual words by clustering using k-means.

    [input]
    * opts          : options
    * n_worker      : number of workers to process in parallel

    [saved]
    * dictionary   : numpy.ndarray of shape (K,3F)
    """

    data_dir = opts.data_dir
    feat_dir = opts.feat_dir
    out_dir = opts.out_dir
    K = opts.K

    train_files = open(join(data_dir, "train_files.txt")).read().splitlines()

    if not os.path.exists(feat_dir):
        os.makedirs(feat_dir)

    if not os.path.exists(out_dir):
        os.makedirs(out_dir)
```

Q1.2

```
# Copy and paste your code here.

mp_inputs = [
    (
        opts,
        join(data_dir, img_path),
        join(feat_dir, "__".join(os.path.splitext(img_path)[0].split("/")) +
             ".npy"),
    )
    for img_path in train_files
]

feature_files = []
with multiprocessing.Pool(n_worker) as pool:
    for result in pool imap_unordered(compute_dictionary_one_image,
                                       mp_inputs):
        feature_files.append(result)

assert len(feature_files) == len(mp_inputs), print(
    "Not all images processed while creating BOW"
)

features = []
for feature_file in feature_files:
    features.append(np.load(feature_file))

features = np.concatenate(features, axis=0)

# example code snippet to save the dictionary
kmeans = KMeans(n_clusters=K).fit(features)
dictionary = kmeans.cluster_centers_

np.save(join(out_dir, "dictionary.npy"), dictionary)
```

1.3 Computing Visual Words

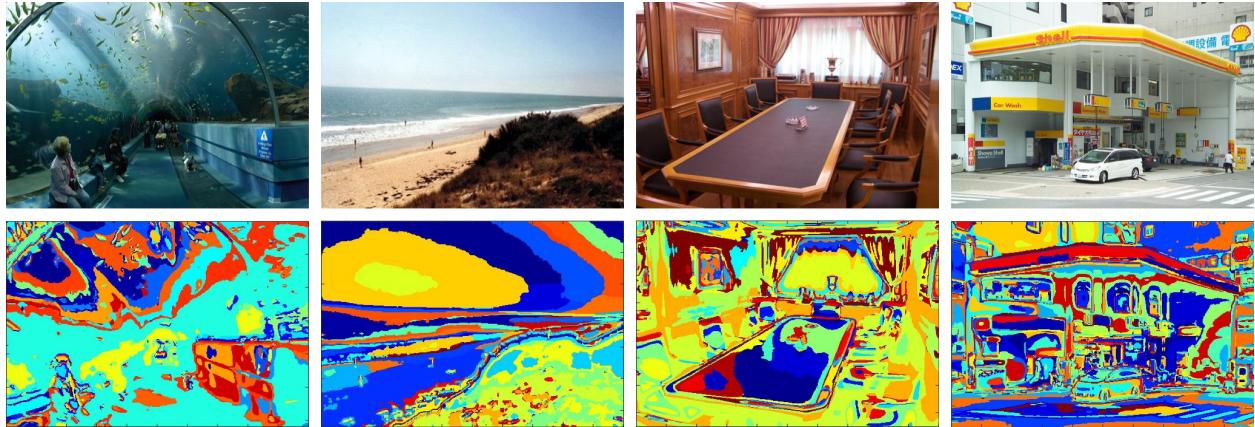


Figure 1.4: Visual words over images. You will use the spatially unordered distribution of visual words in a region (a bag of visual words) as a feature for scene classification, with some coarse information provided by spatial pyramid matching [4].

Q1.3 (10 points): We want to map each pixel in the image to its closest word in the dictionary. Complete the following function to do this:

```
visual_words.get_visual_words(opts, img, dictionary)
```

and return `wordmap`, a matrix with the same width and height as `img`, where each pixel in `wordmap` is assigned the closest visual word of the filter response at the respective pixel in `img`. We will use the standard Euclidean distance to do this; to do this efficiently, use the function `scipy.spatial.distance.cdist()`. Some sample results are shown in Fig 2.2.

Visualize wordmaps for three images. Include some comments on these visualizations: do the “word” boundaries make sense to you? The visualizations should look similar to the ones in Fig 2.2. Don’t worry if the colors don’t look the same, newer `matplotlib` might use a different color map.

We provide the skeleton code from line 28-33 in `main.py`. You can get the results by running `python main.py --filter-scales 1 2 4 --feat-dir TMP_OUT_DIR_FOR_EACH_IMG --out-dir FINAL_OUT_DIR`.

Q1.3 Solution

The word boundaries does make sense as we can see from the wordmaps. Visually similar features have been grouped into it's own cluster. Although for kitchen image, due to having more distinctive features (e.g. sharp edges), the wordmap is much better than the wordmaps obtained for the windmill and aquarium image. For the aquarium image, the people have been grouped together but the surroundings have almost merged together. Probably, increasing the number of clusters would be beneficial to solve this issue.

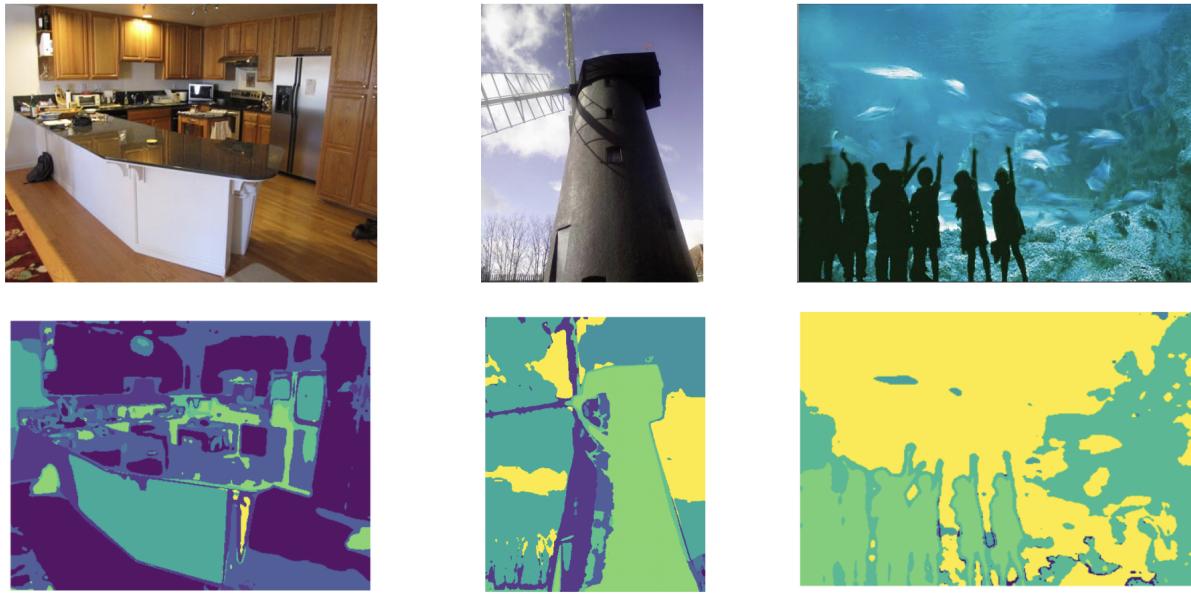


Figure 1.5: Visual words over 3 images using information provided by spatial pyramid matching [4].

2 Building a Recognition System

We have formed a convenient representation for recognition. We will now produce a basic recognition system with spatial pyramid matching. The goal of the system is presented in Fig 0.1: given an image, classify (i.e., recognize/name) the scene depicted in the image.

Traditional classification problems follow two phases: training and testing. At training time, the computer is given a pile of formatted data (*i.e.*, a collection of feature vectors) with corresponding labels (*e.g.*, “desert”, “park”) and then builds a model of how the data relates to the labels (*e.g.*, “if green, then park”). At test time, the computer takes features and uses these rules to infer the label (*e.g.*, “this is green, therefore it is a park”).

In this assignment, we will use the simplest classification method: nearest neighbor. At test time, we will simply look at the query’s nearest neighbor in the training set and transfer that label. In this example, you will be looking at the query image and looking up its nearest neighbor in a collection of training images whose labels are already known. This approach works surprisingly well given a huge amount of data. (For a cool application, see the work by Hays & Efros [3]).

The key components of any nearest-neighbor system are:

- features (how do you represent your instances?) and
- similarity (how do you compare instances in the feature space?).

You will implement both.

2.1 Extracting Features

We will first represent an image with a bag of words. In each image, we simply look at how often each word appears.

Q2.1 (10 points): Write the function

```
visual_recog.get_feature_from_wordmap(opts, wordmap)
```

that extracts the histogram (`numpy.histogram()`) of visual words within the given image (*i.e.*, the bag of visual words). As output, the function will return `hist`, an “ L_1 normalized” dict-size-length histogram. The L_1 normalization makes the sum of the histogram equal to 1. You may wish to load a single visual word map, visualize it, and verify that your function is working correctly before proceeding.

Include your implemented functions within the `minted` block below.

Q2.1

```
# Copy and paste your code here.
def get_feature_from_wordmap(opts, wordmap):
    """
    Compute histogram of visual words.

    [input]
    * opts      : options
    * wordmap   : numpy.ndarray of shape (H,W)

    [output]
    * hist: numpy.ndarray of shape (K)
    """
    K = opts.K
    bins = np.arange(0, K + 1)
    histogram = np.histogram(wordmap, bins=bins, density=True)[0]
    return histogram
```

2.2 Multi-resolution: Spatial Pyramid Matching

A bag of words is simple and efficient, but it discards information about the spatial structure of the image and this information is often valuable. One way to alleviate this issue is to use spatial pyramid matching [4]. The general idea is to divide the image into a small number of cells, and concatenate the histogram of each of these cells to the histogram of the original image, with a suitable weight.

Here we will implement a popular scheme that chops the image into $2^l \times 2^l$ cells where l is the layer number. We treat each cell as a small image and count how often each visual word appears. This results in a histogram for every single cell in every layer. Finally to represent the entire image, we concatenate all the histograms together after normalization by the total number of features in the image. If there are $L + 1$ layers and K visual words, the resulting vector has dimension $K \sum_{l=0}^L 4^l = K (4^{(L+1)} - 1) / 3$.

Now comes the weighting scheme. Note that when concatenating all the histograms, histograms from different levels are assigned different weights. Typically (and in the original work [4]), a histogram from layer l gets half the weight of a histogram from layer $l + 1$, with the exception of layer 0, which is assigned a weight equal to layer 1. A popular choice is to set the weight of layers 0 and 1 to 2^{-L} , and set the rest of the weights to 2^{l-L-1} (e.g., in a three layer spatial pyramid, $L = 2$ and weights are set to $1/4$, $1/4$ and $1/2$ for layer 0, 1 and 2 respectively. See Fig 2.1 for an illustration of a spatial pyramid. Note that the L_1 norm (absolute values of all dimensions summed up together) for the final vector is 1.

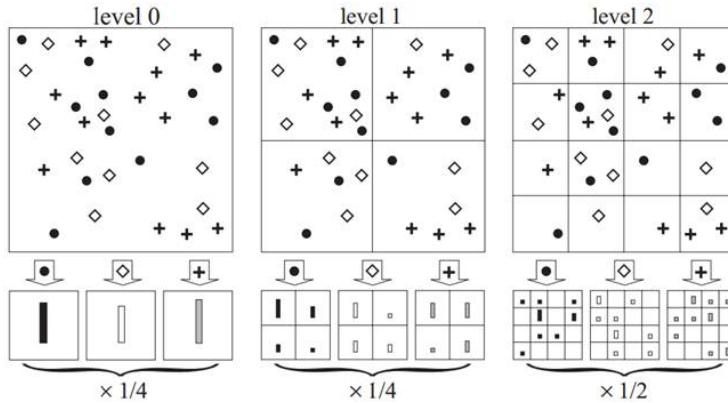


Figure 2.1: Spatial Pyramid Matching: From [4]. Toy example of a pyramid for $L = 2$. The image has three visual words, indicated by circles, diamonds, and crosses. We subdivide the image at three different levels of resolution. For each level of resolution and each channel, we count the features that fall in each spatial bin. Finally, weight each spatial histogram.

Q2.2 (15 points): Create a function `get_feature_from_wordmap_SPM` that forms a multi-resolution representation of the given image.

```
visual_recog.get_feature_from_wordmap_SPM(opts, wordmap)
```

You need to specify the layers of pyramid in `opts.L` (Note there are $L + 1$ layers in total). As output, the function will return `hist_all`, a vector that is L_1 normalized.

One small hint for efficiency: a lot of computation can be saved if you first compute the histograms of the *finest* layer, because the histograms of coarser layers can then be aggregated from finer ones. Make sure you normalize the histogram after aggregation.

Include your implemented functions within the `minted` block below.

Q2.2

```
# Copy and paste your code here.
def get_feature_from_wordmap_SPM(opts, wordmap):
    """
    Compute histogram of visual words using spatial pyramid matching.

    [input]
    * opts      : options
    * wordmap   : numpy.ndarray of shape (H,W)

    [output]
    * hist_all: numpy.ndarray of shape K*(4^(L+1) - 1) / 3
    """
    K = opts.K
    L = opts.L
    h, w = wordmap.shape
    blocks = 2**L
    row_step = h // blocks
    col_step = w // blocks
    total_blocks = blocks * blocks
    histograms = []
    hist_matrix = np.zeros((blocks, blocks, K))
    for block_num in range(total_blocks):
        col = block_num % blocks
        row = block_num // blocks
        local_map = wordmap[
            row * row_step : (row * row_step) + row_step,
            col * col_step : (col * col_step) + col_step,
        ]
        # breakpoint()
        hist_matrix[row, col, :] = get_feature_from_wordmap(opts, local_map)
        histograms.append(0.5 * hist_matrix[row, col, :])

    prev = hist_matrix
    prev_blocks = blocks
    for level in range(L - 1, -1, -1):
        weight = math.pow(2, level - L - 1) if level > 1 else math.pow(2, -L)
        curr_blocks = prev_blocks // 2
        curr_row_step = prev_blocks // curr_blocks
        curr_col_step = prev_blocks // curr_blocks
        curr = np.zeros((curr_blocks, curr_blocks, K))
        total_blocks = curr_blocks * curr_blocks
        for block_num in range(total_blocks):
            col = block_num % curr_blocks
            row = block_num // curr_blocks
            curr[row, col, :] = weight * prev[row * curr_row_step : (row * curr_row_step) + curr_row_step, col * curr_col_step : (col * curr_col_step) + curr_col_step, :].sum(2)
```

Q2.2

```
# Copy and paste your code here.

local_map = prev[
    row * curr_row_step : (row * curr_row_step) + curr_row_step,
    col * curr_col_step : (col * curr_col_step) + curr_col_step,
]
# breakpoint()
curr[row, col, :] = np.sum(local_map, axis=(0, 1)) / (
    curr_row_step * curr_col_step
)
histograms.append(weight * curr[row, col, :])
prev_blocks = curr_blocks
prev = curr

# breakpoint()
histograms = np.asarray(histograms).flatten()
return histograms
```

2.3 Comparing images

We need a way to compare images, to find the “nearest” instance in the training data. In this assignment, we’ll use the histogram intersection similarity. The histogram intersection similarity between two histograms is the sum of the minimum value of each corresponding bins. This is a similarity score: the *largest* value indicates the “nearest” instance.

Q2.3 (10 points): Create the function

```
visual_recog.distance_to_set(word_hist, histograms)
```

where `word_hist` is a $K(4^{L+1} - 1)/3$ vector and `histograms` is a $T \times K(4^{L+1} - 1)/3$ matrix containing T features from T training samples concatenated along the rows. This function computes the histogram intersection similarity between `word_hist` and each training sample as a vector of length T and returns one minus the above quantity as a distance measure (distance is the inverse of similarity). Since this is called every time you look up a classification, you will want this to be fast! (Doing a for-loop over tens of thousands of histograms is a bad idea.) Note: `laundromat/sun_afrrjykuhhlwiwun.jpg` has 4 channels instead of 3. Discard the last channel.

Include your implemented functions within the `minted` block below.

Q2.3

```
# Copy and paste your code here.
def distance_to_set(word_hist, histograms):
    """
    Compute distance between a histogram of visual words with all training image
    histograms.

    [input]
    * word_hist: numpy.ndarray of shape (K)
    * histograms: numpy.ndarray of shape (N, K)

    [output]
    * sim: numpy.ndarray of shape (N)
    """
    minima = np.minimum(word_hist, histograms)
    sim = np.sum(minima, axis=1)
    return sim
```

2.4 Building A Model of the Visual World

Now that we've obtained a representation for each image, and defined a similarity measure to compare two spatial pyramids, we want to put everything up to now together.

Simple I/O code has been provided in the respective functions, which include loading the training images specified in `data/train_files.txt` and the filter bank and visual word dictionary from `dictionary.npy`, and also saving the learned model to `trained_system.npz`. Specifically in `trained_system.npz`, you should have:

1. `dictionary`: your visual word dictionary.
2. `features`: an $N \times K (4^{(L+1)} - 1) / 3$ matrix containing all of the histograms of the N training images in the data set.
3. `labels`: an N vector containing the labels of each of training images. (`features[i]` will correspond to label `labels[i]`).
4. `SPM_layer_num`: the number of spatial pyramid layers you used to extract the features for the training images.

Do not use the testing images for training!

The table below lists the class names that correspond to the label indices:

0	1	2	3	4	5	6	7
aquarium	desert	highway	kitchen	laundromat	park	waterfall	windmill

Q2.4 (15 points): Implement the function

```
visual_recog.build_recognition_system()
```

that produces `trained_system.npz`. You may include any helper functions you write in `visual_recog.py`.

Implement

```
visual_recog.get_image_feature(opts, img_path, dictionary)
```

that loads an image, extract word map from the image, computes the SPM, and returns the computed feature. Use this function in your `visual_recog.build_recognition_system()`.

We provide the skeleton code from line 36-37 in `main.py`. You can train the model by running `python main.py --filter-scales 1 2 4 --feat-dir TMP_OUT_DIR_FOR_EACH_IMG --out-dir FINAL_OUT_DIR`.

Include your implemented functions within the `minted` block below.

Q2.4

```
# Copy and paste your code here.
def get_image_feature(args):
    """
    Extracts the spatial pyramid matching feature.

    [input]
    * arg.opts      : options
    * arg.img_path  : path of image file to read
    * arg.dictionary: numpy.ndarray of shape (K, 3F)

    [output]
    * feature: numpy.ndarray of shape (K*(4^(L+1) - 1) / 3)
    """
    opts, img_path, dictionary = args
    with Image.open(img_path) as img:
        img = np.array(img).astype(np.float32) / 255
        wordmap = visual_words.get_visual_words(opts, img, dictionary)
        spm_feats = get_feature_from_wordmap_SPM(opts, wordmap)
    return spm_feats

def build_recognition_system(opts, n_worker=1):
    """
    Creates a trained recognition system by generating training features from all
    ↵ training images.

    [input]
    * opts      : options
    * n_worker  : number of workers to process in parallel

    [saved]
    * features: numpy.ndarray of shape (N,M)
    * labels: numpy.ndarray of shape (N)
    * dictionary: numpy.ndarray of shape (K,3F)
    * SPM_layer_num: number of spatial pyramid layers
    """
    data_dir = opts.data_dir
    out_dir = opts.out_dir
    SPM_layer_num = opts.L

    train_files = open(join(data_dir, "train_files.txt")).read().splitlines()
    train_labels = np.loadtxt(join(data_dir, "train_labels.txt"), np.int32)
    dictionary = np.load(join(out_dir, "dictionary.npy"))
```

Q2.4

```
# Copy and paste your code here.
mp_inputs = [
    (opts, join(data_dir, img_path), dictionary) for img_path in train_files
]
features = []
with multiprocessing.Pool(n_worker) as pool:
    for result in pool imap_unordered(get_image_feature, mp_inputs):
        features.append(result)

features = np.stack(features, axis=0)

# example code snippet to save the learned system
np.savez_compressed(
    join(out_dir, "trained_system.npz"),
    features=features,
    labels=train_labels,
    dictionary=dictionary,
    SPM_layer_num=SPM_layer_num,
)
```

2.5 Quantitative Evaluation

Qualitative evaluation is all well and good (and very important for diagnosing performance gains and losses), but we want some hard numbers.

Load the test images and their labels, and compute the predicted label of each one. That is, compute the test image's distance to every image in the training set, and return the label of the closest training image. To quantify the accuracy, compute a confusion matrix C . In a classification problem, the entry $C(i, j)$ of a confusion matrix counts the number of instances of class i that were predicted as class j . When things are going well, the elements on the diagonal of C are large, and the off-diagonal elements are small. Since there are 8 classes, C will be 8×8 . The accuracy, or percent of correctly classified images, is given by the trace of C divided by the sum of C . **Hint:** The accuracy with default parameters is 50%.

Q2.5 (10 points): Implement the function

```
visual_recog.evaluate_recognition_system()
```

that tests the system and outputs the confusion matrix. **Report the (a) confusion matrix and your (b) overall accuracy.** This does not have to be formatted prettily: *e.g.*, you can simply copy/paste it into a verbatim environment.

Q2.5 (a) (b)

```
(cv_hw_1) bevanj@Balasubramanyams-MacBook-Air code % python main.py --filter-scales 1 2 4 --feat-dir ../feat_dir --out-dir ../other_out_dir
Params: K = 10 | alpha = 25 | filter_scales = [1.0, 2.0, 4.0] | L = 1
[[32 1 1 4 1 2 4 5]
 [ 0 27 8 5 5 0 2 3]
 [ 1 6 30 0 0 2 2 9]
 [ 3 2 3 31 8 2 0 1]
 [ 4 2 0 14 20 4 3 3]
 [ 3 1 4 3 4 29 3 3]
 [ 3 0 2 2 8 12 20 3]
 [ 5 7 6 1 1 6 3 21]]
0.525
```

Using default parameters, was able to achieve the above confusion matrix. We can see that aquarium and desert are classified the best. And, there are considerable misclassifications for waterfall being misclassified as park and Laundromat being classified as Kitchen.

Overall accuracy was 52.5% using default parameters.

2.6 Find the failures

There are some classes/samples that are more difficult to classify than the rest using the bags-of-words approach. As a result, they are classified incorrectly into other categories.

Q2.6 (5 points): Include some images of these hard classes/samples, and discuss why they are more difficult than the rest.

Q2.6

Misclassification seems to happen due to the presence of common features between classes.

Based on the confusion matrix, misclassified samples seems to come from following places:

1. Laundromat samples being classified as Kitchen
2. Waterfall samples being classified as Park
3. Desert samples being classified as Highway
4. Windmill samples being classified as Highway

Please refer to the images in the next page

1. For the first row, the laundromat image classified is classified as kitchen by the system (using default parameters) due to the fact that in the train samples of laundromat there are some outliers as shown which actually seem to be images of Kitchen. Along with this, the fact that both classes seem to have similar equipments makes the classification task tough
2. For the second row, as shown the waterfall image (left) has been classified as Park due. This seems to be due to the present of water and green surface (trees etc.) as shown in a sample image from park category from train which is present in both parks and waterfall.
3. The third row shows the image of windmill sample (left) which was classified as highway. This probably due to the fact that both these classes seem to have roads which might confuse the recognition system.
4. The final row shows an image of desert sample which was classified as Highway. As can be seen from right image in the final row, samples from Highway do contain flat surfaces with similar texture which might confuse the recognition system



laundromat/sun_alhxqkbvwmxvgdzx.jpg



laundromat/sun_aiyluzcowlbwxmdb.jpg



laundromat/sun_azdzvpzrbrqcmsgssr.jpg



waterfall/sun_awfkpoifyaxauwswo.jpg



park/sun_adfwbwivvtfnwhwd.jpg



windmill/sun_bzepllyllrvbae.jpg



highway/sun_bejfepafwmousxam.jpg



desert/sun_bcmhdbzhuexunnag.jpg



highway/sun_bivkptxtgrmzlbby.jpg

Figure 2.2: Misclassified samples, Left image from test set and Right Image from train set

3 Improving performance

3.1 Hyperparameter tuning

Now we have a full-fledged recognition system plus an evaluation system, it's time to boost up the performance. In practice, it is most likely that a model will not work well out-of-the-box. It is important to know how to tune a visual recognition system for the task at hand.

Q3.1 (15 points): Tune the system you build to reach around 65% accuracy on the provided test set (data/test_files.txt). A list of hyperparameters you should tune is provided below. They can all be found in opts.py.

- filter_scales: a list of filter scales used in extracting filter response;
- K: the number of visual words and also the size of the dictionary;
- alpha: the number of sampled pixels in each image when creating the dictionary;
- L: the number of spatial pyramid layers used in feature extraction.

(a) Include a table of ablation study containing at least 3 major steps (changing parameter X to Y achieves accuracy Z%). (b) Also, describe why you think changing a particular parameter should increase or decrease the overall performance in the table you show.

Q3.1 (a) (b)

By using the default parameters we were able to achieve an accuracy of 52.5%. In order to perform an ablation study to determine affect of hyperparameters, one hyperparameter is changed at a time in order to note the affect of each.

Ablation Study					
filter_scales	K	L	alpha	accuracy	
Change in filter scales					
[1,2,4]	10	1	25	0.525	comparision with default
[1,2,4,8]	10	1	25	52	decreased 0.5%
[1,2,4,8, 11]	10	1	25	0.535	increased, 1%
Change in L					
[1,2,4]	10	1	25	0.525	comparision with default
[1,2,4]	10	2	25	0.5225	decreased 0.5%
[1,2,4]	10	3	25	0.54	increased 1.5%
Change in alpha					
[1,2,4]	10	1	25	0.525	comparision with default
[1,2,4]	10	1	50	0.5225	decreased 0.25%
[1,2,4]	10	1	100	0.52	decreased 0.5%
Change in K					
[1,2,4]	10	1	25	0.525	comparision with default
[1,2,4]	30	1	25	0.61	increased 8.5%
[1,2,4]	50	1	25	0.59	increased 6.5%
[1,2,4]	100	1	25	0.605	increased 8%

Based on the above table we can draw the following conclusions:

- Change in K, the number of clusters really increased the accuracy of the recognition system. This is something that can be expected since as we increase the number of clusters, the average distortion should decrease as well
- When multiple scales are used we expect to accumulate different scales of features as well which should help in the robustness of the recogniton system. And as seen from the table ([1,2,4] to [1,2,4,8,11]), we do see an increase

Q3.1 (a) (b)

- When alpha is increased, we notice decrease in overall accuracy. But, I believe this result needs to be interpreted with caution since alpha varies the number of pixels we are considering in the image. And more number of pixels (more data) is mostly always a good thing. Probably there needs to be a change in filter scales to draw out meaningful features
- When L is increased from 1 to 3 we see an increase in performance. This can be attributed to the fact that by increasing the spatial layers we expect to capture more details which would be useful

The confusion matrix and the accuracy of the best performing model which crossed the set threshold of 60% is as below. Note: The uploaded code contains the best model which obtained better accuracy than this.

```
(cv_hw_1) bevani@Balasubramanyams-MacBook-Air code % python main.py --feat-dir ..//hyp_feat_dir --out-dir ..//hyp_out_dir
> Params: K = 30 | alpha = 25 | filter_scales = [1, 2, 4] | L = 1
[[36  0   0   1   5   1   2   5]
 [2 27  5   8   2   1   2   3]
 [1  4 30  0   1   4   3   7]
 [1  0  1 39  8   1   0   0]
 [0  0  0 15 22  7   6   0]
 [4  0  2  2  2 34  2   4]
 [4  1  3  1  5  8 27  1]
 [2  4  8  0  1  3  3 29]]
0.61
```

3.2 [Extra Credit] Further improvement

Q3.2 (10 points): Can you improve your classifier, in terms of accuracy or speed? Be creative! Or be well-informed, and cite your sources! For some quick ideas, try resizing the images, subtracting the mean color, changing the structure or weights of the spatial pyramid, or replacing the histogram intersection with some other similarity score. Whatever you do, explain:

- (1) what you did,
- (2) what you expected would happen, and
- (3) what actually happened.

Include these results in the report and submit the code.

Q3.2

I believe I was able to implement multiprocessing wherever necessary to improve the speed. My local system had 8 cpu cores and hence took advantage of it to train and evaluate the recognition system

- Used multiprocessing to compute dictionary for faster processing (ref. visual_words.py, compute_dictionary function)
- Used multiprocessing to compute image features while training and evaluating (ref. visual_recog.py, build_recognition_system, evaluate_recognition_system functions)
- Also used multiprocessing to evaluate each image during evaluation phase (ref. visual_recog.py, evaluate_recognition_system function)

In terms of increasing the accuracy based on the previous question's observations I opted to increase the number of clusters K to 100, L to 3 and opted for filter scales [1, 2, 4, 8]. For the alpha value, I used 100, although I observed decreased performance in my ablation I believe it's highly dependent on L and how the pixels were sampled.

I expected the accuracy to increase from the previous version and in line with my expectations based on the aforementioned hyperparameters, **I was able to increase the accuracy to 65%**. Though the interesting thing is the variance, due to random sampling of pixels, I have observed performance accuracy range from 64.5% to as high as 65%. Though it varies, the variance is low, hence chose to stick with this model. The confusion matrix is as below and the dictionary and model have been uploaded along with the code.

```
Namespace(K=100, L=array(3), alpha=100, data_dir='..../data', feat_dir='..../optimal_feat_dir', filter_scales=[1, 2, 4, 8], out_dir='..../o
ptimal_out_dir')
[[42  0   0   1   3   0   2   2]
 [ 0 36  2   2   3   3   0   4]
 [ 2  4 30  0   0   4   0 10]
 [ 1  3  1 31 11  2   0   1]
 [ 0  2  2  5 26  9   3   3]
 [ 2  0  2  1  3 34  3   5]
 [ 4  0  5  0  3  5 30  3]
 [ 1  4 10  0  0  3  1 31]]
0.65
```

4 HW1 Distribution Checklist

After unpacking `hw1.zip`, you should have a folder `hw1` containing one folder for the data (`data`), one for your code (`code`), and one for the report (`latex`). In the `code` folder, where you will primarily work, you will find:

- `visual_words.py`: function definitions for extracting visual words.
- `visual_recog.py`: function definitions for building a visual recognition system.
- `util.py`: some utility functions
- `main.py`: main function for running the system

The data folder contains:

- `data/`: a directory containing `.jpg` images from the SUN database.
- `data/train_files.txt`: a text file containing a list of training images.
- `data/train_labels.txt`: a text file containing a list of training labels.
- `data/test_files.txt`: a text file containing a list of testing images.
- `data/test_labels.txt`: a text file containing a list of testing labels.

5 HW1 submission checklist

Submit your write-up and code to Gradescope.

- **Writeup.** Please use this provided template for your writeup. The write-up should be a pdf file named `<AndrewId>.hw1.pdf`. **You must select the pages of the writeup that correspond to each question.**
- **Code.** The code should be submitted as a zip file named `<AndrewId>.hw1.zip`. By extracting the zip file, it should have the following files in the structure defined below.

When you submit, remove the folder `data/` and `feat/` if applicable, as well as any large temporary files that we did not ask you to create.

- `<andrew.id>/` # A directory inside .zip file
 - * `code/`
 - `dictionary.npy`
 - `trained_system.npz`
 - `<!– all of your .py files >`
 - * `<andrew.id>.hw1.pdf` make sure you upload this pdf file to Gradescope. Please assign the locations of answers to each question on Gradescope.

References

- [1] K. Chatfield, V. Lempitsky, A. Vedaldi, and A. Zisserman. The devil is in the details: an evaluation of recent feature encoding methods. In *British Machine Vision Conference*, 2011.

- [2] K. Grauman and T. Darrell. The pyramid match kernel: discriminative classification with sets of image features. In *Computer Vision (ICCV), 2005 IEEE International Conference on*, volume 2, pages 1458–1465 Vol. 2, 2005.
- [3] James Hays and Alexei A Efros. Scene completion using millions of photographs. *ACM Transactions on Graphics (SIGGRAPH 2007)*, 26(3), 2007.
- [4] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *Computer Vision and Pattern Recognition (CVPR), 2006 IEEE Conference on*, volume 2, pages 2169–2178, 2006.
- [5] D.G. Lowe. Object recognition from local scale-invariant features. In *Computer Vision (ICCV), 1999 IEEE International Conference on*, volume 2, pages 1150–1157 vol.2, 1999.
- [6] Laura Walker Renninger and Jitendra Malik. When is scene identification just texture recognition? *Vision research*, 44(19):2301–2311, 2004.
- [7] J. Winn, A. Criminisi, and T. Minka. Object categorization by learned universal visual dictionary. In *Computer Vision (ICCV), 2005 IEEE International Conference on*, volume 2, pages 1800–1807 Vol. 2, 2005.
- [8] Jianxiong Xiao, J. Hays, K.A. Ehinger, A. Oliva, and A. Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 3485–3492, 2010.