

CS1674: Homework 4

Due: 3/2/2021, 11:59pm

This assignment is worth 50 points.

Download HW4_starter.zip from the Canvas HW4 Submission Page. (similar to how you got hw2_starter.zip for HW2)

Part I: Feature Description (30 points)

In this problem, you will implement a feature description pipeline, as discussed in class. While you will not exactly implement that, [the SIFT paper](#) by David Lowe is a useful resource, in addition to Section 4.1 of the Szeliski textbook.

Use the following signature: `function [features] = compute_features(x, y, scores, Ix, Iy)` (note $I_x=I_h$, $I_y=I_v$). The inputs x , y , $scores$, I_x , I_y are defined in HW3. The output `features` is an $N \times d$ matrix, each row of which contains the d -dimensional descriptor for the n -th keypoint. We'll simplify the histogram creation procedure a bit, compared to the original implementation presented in class. In particular, we'll compute a descriptor with dimensionality $d=8$ (rather than $4 \times 4 \times 8$), which contains an 8-dimensional histogram of gradients computed from a 11×11 grid centered around each detected keypoint (i.e. $-5:+5$ neighborhood horizontally and vertically).

1. [4 pts] If any of your detected keypoints are less than 5 pixels from the top/left or 5 pixels from the bottom/right of the image, i.e. pixels lacking $5+5$ neighbors in either the horizontal or vertical direction, erase this keypoint from the x , y , $scores$ vectors *at the start of your code* and do not compute a descriptor for it.
2. [8 pts] To compute the gradient magnitude $m(x, y)$ and gradient angle $\theta(x, y)$ at point (x, y) , take I to be the image and use the formula below shown in class and Matlab's `atan2`, which returns values in the range $[-90, 90]$. If the gradient magnitude is 0, then both the x and y gradients are 0, and you should ignore the orientation for that pixel (since it won't contribute to the histogram).

$$m(x, y) = \sqrt{I_x^2(x, y) + I_y^2(x, y)}$$
$$\theta(x, y) = \tan^{-1}(I_y(x, y)/I_x(x, y))$$

Note that I_x , I_y can be derived the usual way ('prewitt', 'sobel', etc.) following however you computed them in `extract_keypoints` from HW3.

3. [6 pts] Quantize the gradient orientations in 8 bins (so put values between -90 and -67.5 degrees in one bin, the -67.5 to -45 degree angles in another bin, etc.) For example, you can have a quantization matrix with the same size as the image, that says to which bin (1 through 8) the gradient of each pixel belongs.
4. [6 pts] Now, let's compute the SIFT histogram for each keypoint (based on its 11×11 grid centered around each detected keypoint, i.e. $-5:+5$ neighborhood horizontally and vertically). To populate the SIFT histogram, consider each of the 8 bins. To populate the first bin, sum the gradient magnitudes that are between -90 and -67.5 degrees. Repeat analogously for all bins.
5. [6 pts] Finally, you should clip all values to 0.2 as discussed in class, and normalize each descriptor to be of unit length, e.g. using `hist_final = hist_final / sum(hist_final);` Normalize both before and after the clipping (i.e., 1. normalize, 2. clip, 3. normalize again). You do not have to implement any more sophisticated detail from the Lowe paper.

Part II: Image Description with SIFT Bag-of-Words (10 points)

In this part, you will compute a bag-of-words histogram representation of an image. Conceptually, the histogram for image I_j is a k -dimensional vector: $F(I_j) = [freq_{1,j} \quad freq_{2,j} \quad \dots \quad freq_{k,j}]$, where each entry $freq_{i,j}$ counts the number of occurrences of the i -th visual word in image j , and k is the number of total words in the vocabulary. (Acknowledgement: Notation from Kristen Grauman's assignment.)

Use the following function signature: `function [bow_repr] = computeBOWRepr(features, means)` where `bow_repr` is a normalized bag-of-words histogram (a *row* vector), `features` is the $N \times 8$ set of descriptors computed for the image (output by the function you implemented in Part I above), and `means` is a $k \times 8$ set of cluster means. A set of files with cluster means are provided for you on Canvas (for $k=\{2, 5, 10, 50, 100, 200\}$).

1. [2 pt] A bag-of-words histogram has as many dimensions as the number of clusters k , so initialize the `bow` variable accordingly.
2. [4 pts] Next, for each feature (i.e. each row in `features`), compute its distance to each of the cluster means, and find the closest mean. A feature is thus conceptually "mapped" to the closest cluster. You can do this efficiently using Matlab's `pdist2` function (with inputs `features`, `means`).
3. [4 pts] To compute the bag-of-words histogram, count how many features are mapped to each cluster.
4. [2 pts] Finally, normalize the histogram by dividing each entry by the sum of the entries. (just a simple normalization. No clipping.)

Part III: Comparison of Image Descriptors (10 points)

In this part, we will test the quality of the different representations. A good representation is one that retains some of the semantics of the image; oftentimes by "semantics" we mean object class label. In other words, a good representation should be one such that two images of the same object have similar representations, and images of different objects have different representations. We will test to what extent this is true, using our images of cardinals, leopards and pandas.

To test the quality of the representations, we will compare two averages: the average *within-class distance* and the average *between-class distance*. A representation is a vector, and "distance" is the Euclidean distance between two vectors (i.e. the representations of two images). "Within-class distances" are distances computed between the vectors for images of the same class (i.e. leopard-leopard, panda-panda). "Between-class distances" are those computed between images of different classes, e.g. leopard-panda. If you have a good image representation, should the average within-class or the average between-class distance be smaller?

A script `compare_representations.m` is provided for you, which does the following:

1. Reads the images, and resizes them to 300x300.
2. Uses the (a) code you wrote above, (b) the function `extract_keypoints` you wrote in HW3 (talk to the TA or instructor if that function doesn't work), and (c) another function `[texture_repr_concat, texture_repr_mean] = computeTextureReprs(image, F)` provided for you, to compute three types image representations (bag of words using SIFT descriptors using different values of k i.e. different number of cluster centers, concatenation of the texture representation of all pixels, and average over the texture representation of pixels) for each image.
3. Computes and prints the ratio *average_within_class_distance / average_between_class_distance* for each representation type.
4. `compare_representations` calls `computeTextureReprs` which has the following inputs:
image is the output of an `imread`, and `F` is the $49 \times 49 \times \text{num_filters}$ matrix of filters you used in HW2.
The latter function computes two image representations based on the filter responses. The first is simply a

concatenation of the filter responses for all pixels and all filters. The second contains the mean of filter responses (averaged across all pixels) for each image.

Your task in this part is to do the following. Run `compare_representations` (check the top of the script for use tips); it tries different values of k . Experiment with different values for the number of keypoints you extract from images, by changing your `extract_keypoints` function (no need to submit any changes to it; we'll use our own function). Then, in a file `answers.txt`, answer the following questions.

- For which of the three representations is the within-between ratio smallest?
- Does the answer to this question depend on the value of k that you use?
- Does it depend on the number of keypoints you extract? (Try 500, 1000, 2000, 3000.)
- Which of the three types of descriptors that you used is the best one? How can you tell?
- Is this what you expected? Why or why not?

You can copy-paste ratios from the output of `compare_representations` as needed to support your answers.

Notes:

1. Do not worry if you see NaN for $k=2$.
2. Do not worry if the patterns are not obvious. They may not be strictly increasing or decreasing.
3. You would expect the within-between ratio to be <1 for some k , since that implies the within-class (i.e., images of same animals) distances are smaller than the between-class (i.e., images of different animals). Thus, checking to see if the within-between ratio is <1 for some k would be the least sanity check.
4. The purpose of Part III is to use the provided script to make observations. This is similar to an actual research experiment. Depending on implementation of your `compute_features.m`, `extract_keypoints.m`, and the cluster centers, your answers may differ. We are looking for your descriptive observations based on the questions above. We are **not** looking to see if you provide the correct k , the number of keypoints, etc., because they depend on many factors.

Submission:

- `compute_features.m` (from Part I)
- `computeBOWRepr.m` (from Part II)
- `answers.txt` (from Part III)

Acknowledgement: Adriana Kovashka.