

Recognizing pictures at an exhibition using SIFT

Rahul Choudhury
Biomedical Informatics Department
Stanford University
EE 368 Project Report

Abstract—To aid visitors of art centers in automatically identifying an object of interest, a computer vision algorithm known as SIFT has been implemented using MATLAB. Salient aspects of the algorithm have been discussed in light of image recognition applications. Test methods to validate the algorithm have been detailed out. Performance improvement for SIFT also has been proposed.

1. INTRODUCTION

Museums and art centers present objects of interests, which contain important domain-specific information that visitors are interested in. Art exhibition centers usually provide a guide or a booklet, which visitors manually scan to find out general information about an object of interest. However this manual and passive method of information retrieval is tedious and interrupts visitor's art-viewing experience. The goal of this project is to develop an automated software-based technique, which can recognize paintings on display in an art-center based on snapshots taken with a camera-phone.

2. BACKGROUND

Comparing images in order to establish a degree of similarity is an important computer vision problem and has application in various domains such as interactive museum guide, robot localization, content-based medical image retrieval, and image registration. Comparing images remains a challenging task because of issues such as variation in illumination conditions, partial occlusion of objects, differences in image orientation etc. Global image characteristics such as color histograms, responses to filter banks etc. are usually not effective in solving real-life image matching problems.

Researchers have recently turned their attention to local features in an image, which are invariant to common image transformations and variations. Usually two broad steps are found in any local feature-based image-matching scheme. The first step involves detecting features (also referred to as keypoints or interest points) in an image in a repeatable way. Repeatability is important in this step because robust

matching cannot be performed if the detected locations of keypoints on an object vary from image to image. The second step involves computing descriptors for each detected interest point. These descriptors are useful to distinguish between two keypoints.

The goal is to design a highly distinctive descriptor for each interest point to facilitate meaningful matches, while simultaneously ensuring that a given interest point will have the same descriptor regardless of the object's pose, the lighting in the environment, etc. Thus both detection and description steps rely on invariance of various properties for effective image matching.

Image matching techniques based on local features are not new in the computer vision field. Van Gool [5] introduced the generalized color moments to describe the shape and the intensities of different color channels in a local region. Sven Siggelkow [2] used feature histograms for content-based image retrieval. These methods have achieved relative success with 2D object extraction and image matching. Mikolajczyk and Schmid [3] used the differential descriptors to approximate a point neighborhood for image matching and retrieval. Schaffalitzky and Zisserman [4] used Euclidean distance between orthogonal complex filters to provide a lower bound on the Squared Sum Differences (SSD) between corresponding image patches. David Lowe proposed [1] Scale Invariant Feature Transforms (SIFT), which are robustly resilient to several common image transforms. Mikolajczyk and Schmid [6] reported an experimental evaluation of several different descriptors where they found that the SIFT descriptors obtain the best matching results.

3. SELECTION OF LOCAL FEATURES

The following requirements were key in selecting a suitable local-feature for images used in this project:

- (a) Invariance: The feature should be resilient to changes in illumination, image noise, uniform scaling, rotation, and minor changes in viewing direction
- (b) Highly Distinctive: The feature should allow for correct object identification with low probability of

mismatch.

- (c) Performance: Given the nature of the image recognition problem for an art center, it should be relatively easy and fast to extract the features and compare them against a large database of local features.

4. BRIEF DESCRIPTION OF THE ALGORITHM

Based on the requirements mentioned in the previous section and the reported robustness in [6], I selected the Scale Invariant Feature Transform (SIFT) [1] approach for this project. SIFT is an approach for detecting and extracting local feature descriptors that are reasonably invariant to changes in illumination, image noise, rotation, scaling, and small changes in viewpoint. Before performing any multi-resolution transformation via SIFT, the image is first converted to grayscale color representation.

A complete description of SIFT can be found in [1]. An overview of the algorithm is presented here. The algorithm has four major stages as mentioned below:

- Scale-space extrema detection: The first stage searches over scale space using a Difference of Gaussian function to identify potential interest points.
- Keypoint localization: The location and scale of each candidate point is determined and keypoints are selected based on measures of stability.
- Orientation assignment: One or more orientations are assigned to each keypoint based on local image gradients.
- Keypoint descriptor: A descriptor is generated for each keypoint from local image gradients information at the scale found in the second stage.

Each one of the above-mentioned stages is elaborated further in the following sections.

A. Scale-space Extrema Detection

The SIFT feature algorithm is based upon finding locations (called keypoints) within the scale space of an image which can be reliably extracted. These keypoints correspond to local extrema of difference-of-Gaussian (DoG) filters at different scales. The DoG image is represented as $D(x, y, \sigma)$ and can be computed from the difference of two nearby scaled images separated by a multiplicative factor k :

$$\begin{aligned} D(x, y, \sigma) &= (G(x, y, k\sigma) - G(x, y, \sigma)) * I(x, y) \\ &= L(x, y, k\sigma) - L(x, y, \sigma) \end{aligned}$$

where $L(x, y, \sigma)$ is the scale space of an image, built by convolving the image $I(x, y)$ with the Gaussian kernel $G(x, y, \sigma)$. The convolved images are grouped by octave. An octave corresponds to doubling the value of σ , and the value of k is selected so that a fixed number of blurred images are generated per octave. This also ensures the same number of DoG images are obtained per octave. Keypoints are identified as local maxima or minima of the DoG images across scales. Each pixel in a DoG image is compared to its 8 neighbors at the same scale, and the 9 corresponding neighbors at neighboring scales. If the pixel is a local maximum or minimum, it is selected as a candidate keypoint.

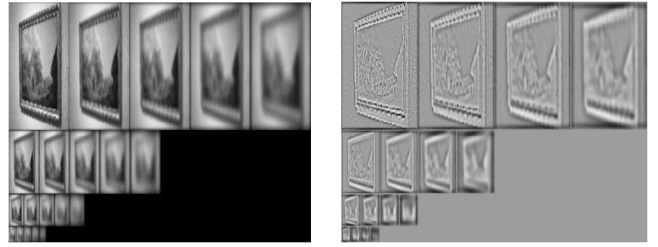


Figure 1: Gauss Pyramid (left) and DoG (right) for a sample training image

B. Keypoint Localization

In this step the keypoints are filtered so that only stable and more localized key points are retained. First a 3D quadratic function is fitted to the local sample points to determine the location of the maximum. If it is found that the extremum lies closer to a different sample point, the sample point is changed and the interpolation performed instead about that point. The function value at the extremum is used for rejecting unstable extrema with low contrast.

The DoG operator has a strong response along edges present in an image, which give rise to unstable key points. A poorly defined peak in the DoG function will have a large principal curvature across the edge but a small principal curvature in the perpendicular direction. The principal curvatures can be computed from a 2x2 Hessian matrix H , computed at the location and scale of the keypoint. H is given by $H = \begin{bmatrix} D_{xx} & D_{xy} \\ D_{yx} & D_{yy} \end{bmatrix}$. The eigen-

values of H are proportional to the principal curvatures of D . However eigenvalues are not explicitly computed; instead trace and determinants of H are used to reject those keypoints for which the ratio between principal curvatures is greater than a threshold.



Figure 2: Keypoints after removing low-contrast interest points



Figure 3: Stable keypoints after removing the edge elements

C. Orientation Assignment

In this step each keypoint is assigned an orientation to make the descriptor invariant to rotation. This keypoint orientation is calculated from an orientation histogram of local gradients from the closest smoothed image $L(x, y, \sigma)$. For each image sample $L(x, y)$ at this scale, the gradient magnitude $m(x, y)$ and orientation $\theta(x, y)$ is computed using pixel differences:

$$m(x, y) = \sqrt{((L(x+1, y) - L(x-1, y)))^2 + ((L(x, y+1) - L(x, y-1)))^2}$$

$$\theta(x, y) = \tan^{-1} \left(\frac{L(x, y+1) - L(x, y-1)}{L(x+1, y) - L(x-1, y)} \right)$$

The orientation histogram has 36 bins covering the 360 degrees of major orientation bins. Each point is added to the histogram weighted by the gradient

magnitude $m(x, y)$ and by a circular Gaussian with σ that is 1.5 times the scale of the keypoint. Additional keypoints are generated for keypoints with multiple local dominant peaks whose magnitude is within 80% of each other. The dominant peaks in the histogram are interpolated with their neighbors for a more accurate orientation assignment.

D. Keypoint Descriptor

The local gradient data from the closest smoothed image $L(x, y, \sigma)$ is also used to create the keypoint descriptor. This gradient information is first rotated to align it with the assigned orientation of the keypoint and then weighted by a Gaussian with σ that is 1.5 times the scale of the keypoint. The weighted data is used to create a nominated number of histograms over a set window around the keypoint. Typical keypoint descriptors use 16 orientation histograms aligned in a 4x4 grid. Each histogram has 8 orientation bins each created over a support window of 4x4 pixels. The resulting feature vectors are 128 elements with a support window of 16x16 scaled pixels.

5. KEYPOINTS MATCHING

Given a test image, each one of its keypoint is compared with keypoints of every image present in the training database. Euclidean distance between each invariant feature descriptor of the test image and each invariant feature descriptor of a database image is computed at first. However two keypoints with the minimum Euclidean distance (the closest neighbors) cannot necessarily be matched because many features from an image may not have any correct match in the training database (either because of background clutter or may be the feature was not detected at all in the training images). Instead the ratio of distance between closest neighbors and distance between the second closest neighbor is computed. If the ratio of distances is greater than 0.6, then that match is rejected. Once the comparison is performed with all test images in the database, the title of the database image with the maximum number of matches is returned as the recognized image.

6. IMPLEMENTATION

The approach described above has been implemented using MATLAB. The implementation has two aspects: training and inference. During the training phase locally invariant features (keypoints, orientation, scales and descriptors) from all training images are retrieved using

the SIFT algorithm and stored in a file named “training-data.m”. The training matlab function is in the script “train.m”.

During inference, the objective is to recognize a test image. A set of local invariant features are retrieved for the test image during the inference phase and compared against the training feature-set using the metric explained in section 5. The title of the closest match is returned as the final output. The inference matlab function is in the script “identify_painting.m”.

An important aspect of the SIFT algorithm is that it generates a large number of features over a broad range of scales and locations. The number of features generated depends primarily on image size and content, as well as algorithm parameters. As the training images for this project have a resolution of 2048x1536, I have down-sampled the images by a factor of 16 to reduce the number of generated keypoints.

7. RESULTS

I first tested the algorithm with leave-one-out cross validation method. The idea is to repeatedly train the algorithm on all but one training examples in the training dataset and then test the algorithm on the held-out example. The resulting errors are then averaged together to obtain an estimate of the generalization error of the model. The generalization error from leave-one-out cross validation method was found out to be zero.

I also performed tests on 30 images of my own that I acquired from the Cantor Arts center using a camera phone (two of them are shown in Figure 4 and Figure 5). The algorithm proves robust to various scales, orientation, and various illumination conditions. If the viewing angle is extreme, I have found that on certain test images, the algorithm fails to recognize the image.



Figure 4: Test Image acquired from the Cantor Arts Center



Figure 5: Test Image acquired from the Cantor Arts Center

The algorithm is somewhat robust against Gaussian noise as well. In Figure 6, an example is shown where Gaussian noise was added twice to a test image using the MATLAB function ‘imnoise’. The algorithm could still successfully identify the image. However if noise is added more than three times using the ‘imnoise’ function, the algorithm fails to recognize the image.

I also tested the algorithm with several test images after adding “salt and pepper” noise to it (an example is shown in Figure 7). I used the MATLAB function ‘imnoise’ for this purpose as well. The algorithm correctly recognized the test images as long as the density is less than 0.05 for the salt and pepper noise.

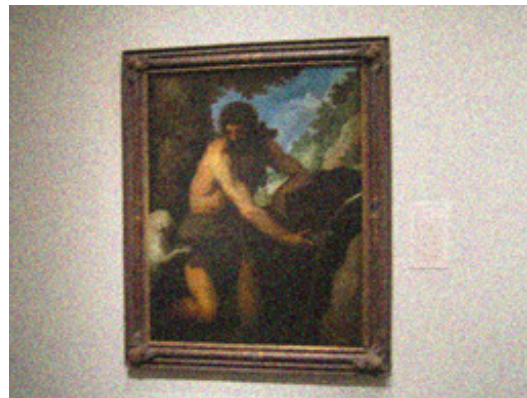


Figure 4: Test image where Gaussian noise was added by using the MATLAB ‘imnoise’ function twice

It takes approximately 40 minutes to train the algorithm in the SCIEN machine using the 99 training images that were provided. Also, for a test image of size 2048x1536, depending on the visual content of the image, it takes anywhere from 15 seconds to 55 seconds for recognition in the SCIEN machine.

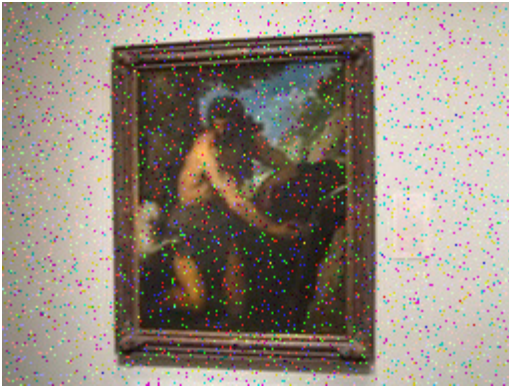


Figure 5: Test image where Salt and Pepper noise was added by using the MATLAB 'imnoise' function

8. FUTURE IMPROVEMENT

In future I would like to work on improving the performance of the SIFT algorithm in several aspects as mentioned below:

A. Global Features

An image may have many areas that are locally similar to each other. These multiple locally similar areas may produce ambiguities when matching local descriptors. The local invariant features of SIFT can be augmented by computing global features of an image.

B. Difference of Mean

As pointed out in [8], instead of using the DoG images for finding the local scale space extrema, we could use Difference of Mean (DoM) images to approximate DoG. To find out DoM images, we would first need to compute an integral image. An integral image can be computed from an input image I as mentioned below:

$$J(x, y) = \sum_{x'=0}^x \sum_{y'=0}^y I(x', y')$$

The mean of a rectangular region of an integral image can be computed very efficiently and independent of the size of the region.

9. CONCLUSION

This project proved to be a compelling exploration of the applications of image processing techniques discussed in EE368. SIFT is a state-of-the-art algorithm for extracting locally invariant features and this project gave me an opportunity to understand several aspects of multi-resolution image analysis and its application in image recognition. I believe this effort resulted in a

robust image recognition implementation, which should perform quite well with the final test images.

ACKNOWLEDGMENT

Author thanks the EE368 teaching staff for their help and encouragement during the quarter.

REFERENCES

- [1] David G. Lowe. Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision*, 60, 2 (2004), pp. 91-110.
- [2] Sven Siggelkow. Feature Histograms for Content-Based Image Retrieval. PhD Thesis, Albert-Ludwigs-University Freiburg, December 2002.
- [3] Mikolajczyk, K., Schmid, C.: An Affine Invariant Interest Point Detector. In: *ECCV*, (2002) 128-142
- [4] Schaffalitzky, F., Zisserman, A.: Multi-view Matching for Unordered Image Sets, or "How Do I Organize My Holiday Snaps?" In: *ECCV*, (2002) 414-431
- [5] Van Gool, T. Moons, and D. Ungureanu. Affine photometric invariants for planar intensity patterns. In *ECCV*, pp. 642-651, 1996.
- [6] Mikolajczyk, K., Schmid, C.: Indexing Based on Scale Invariant Interest Points. In: *ICCV*, (2001) 525-531
- [7] David G. Lowe. Object recognition from local scale-invariant features. *International Conference on Computer Vision*, Corfu, Greece (September 1999), pp. 1150-1157.
- [8] Michael Grabner, Helmut Grabner, and Horst Bischof. Fast approximated SIFT. *Proc. ACCV 2006*, Hyderabad, India