

# Problem Set 4: Harris, SIFT, RANSAC

**Due Monday, March 16, 2015 at 7:00 am EST**

## Description

The focus of this problem set is on feature computation and model fitting. As defined in class features must be (1) fairly repeatable - will tend to show up in both images even with changes in lighting or imaging; (2) well localizable - their location in the imagery should be easily and relatively precisely determined; (3) fairly common without being dense in the imagery; (4) characterizable such that it is possible to find likely matches. Once we have such features and their *putative* matches, we use RANSAC as a form of global checking to find a likely alignment.

You will do each of these steps below, though for some of them code will be provided. SIFT libraries you can use are described in this [supplemental document](#) (see question 2 for instructions).

## What to submit

Download and unzip the ps4 folder (also under: <https://www.udacity.com/wiki/ud810>): [ps4.zip](#)

Rename it to ps4\_xxxx (i.e. ps4\_matlab, ps4\_octave, or ps4\_python) and add in your solutions:

ps4\_xxxx/

- input/ - input images, videos or other data supplied with the problem set
- output/ - directory containing output images and other files your code generates
- ps4.m or ps4.py - code for completing each part, esp. function calls; all functions themselves must be defined in individual function files with filename same as function name, as indicated
- \*.m or \*.py - Matlab/Octave function files (one function per file), Python modules, any utility code
- ps4\_report.pdf - a PDF file with all output images and text responses

Zip it as ps4\_xxxx.zip, and submit on T-Square.

## Guidelines

1. Include all the required images in the report to avoid penalty.
2. Include all the textual responses, outputs and data structure values (if asked) in the report.
3. Make sure you submit the correct (and working) version of the code.
4. Include your name and GTID on the report.
5. Even if the code is not working, submit the code as the instructors can read through the algorithms to give partial credit. Comment your code appropriately.
6. Late submissions should be emailed to the TAs to be graded for partial credit.

## Questions

### 1. Harris corners

In class and in the text we have developed the *Harris* operator. To find the Harris points you need to compute the gradients in both the X and Y directions. These will probably have to be lightly filtered using a Gaussian to be well behaved. You can do this either the “naive” way - filter the image and then do simple difference between left and right (X gradient) or up and down (Y gradient) - or you can take an analytic derivative of a Gaussian in X or Y and use that filter. The scale of the filtering is up to you. You may play with the size of the Gaussian as it will interact with the window size of the corner detection.

a. Write functions to compute both the X and Y gradients. Try your code on both *transA* and *simA*. To display the output, adjoin the two gradient images (X and Y) to make a new, twice as wide, single image (the “gradient-pair”). Since gradients have negative and positive values, you’ll need to produce an image that is gray for 0.0 and black is negative and white is positive.

**Output:** The gradient-pair image for both *transA* and *simA*.

-*transA* gradient-pair image as *ps4-1-a-1.png*

-*simA* gradient-pair image as *ps4-1-a-2.png*

b. Now you can compute the Harris value for the image. As a reminder the Harris value was defined as:

$$R = \det(M) - \alpha \text{trace}(M)^2$$

where

$$M = \sum_{x,y} w(x,y) \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix}$$

The only design decisions are the size of the window (the sum), the windowing function that controls the weights, and the value of  $\alpha$  in the Harris scoring function. Remember you can check you code on the checkerboard images, though those might use a different optimal window size than the real ones.

Write code to compute the Harris value. You can try the weights just equal to 1. But it might work

better with a smoother Gaussian that is higher at the middle and falls off gradually. Your output is a scalar function. Apply to transA, transB, simA, and simB.  
(To display the output reasonably you will have to scale the image values to be in a range of 0-255 or 0.0 to 1.0, depending upon how you deal with images.)

**Output:** The Harris value output image for:

- transA image as ps4-1-b-1.png
- transB image as ps4-1-b-2.png
- simA image as ps4-1-b-3.png
- simB image as ps4-1-b-4.png

c. Finally you can find some corner points. To do this requires two steps: thresholding and non-maximal suppression. You'll need to choose a threshold value that eliminates points that don't seem to be plausible corners. And for the non-maximal suppression, you'll need to choose a radius (could be size of a side of a square window instead of a circle radius) over which a pixel has to be a maximum.

Write a function to threshold and do non-maximal suppression on the Harris output. Surprise, huh? Adjust the threshold and radius until you get a "nice" set of points, probably on the order of a hundred or two (or three?). But use your judgment in terms of getting enough points. Are there any points that are visible in both images but not found as corners in both?

**Output:** Apply your function to both image pairs: (transA, transB) and (simA, simB). Mark the corners visibly in each of the four result images and provide those images.

- transA image with Harris corners marked as ps4-1-c-1.png
- transB image with Harris corners marked as ps4-1-c-2.png
- simA image with Harris corners marked as ps4-1-c-3.png
- simB image with Harris corners marked as ps4-1-c-4.png

**Output (Textual Response):**

- Describe the behavior of your corner detector including anything surprising, such as points not found in both images of a pair.

## 2. SIFT features

Now that you have keypoints for both image pairs, we can compute descriptors. You will be glad to know that we do not expect you to write your own SIFT descriptor code. Instead you'll use a MATLAB package called VLFeat and for Python, the SIFT or SURF classes in OpenCV. In our past experience using VLFeat with Python is not a good idea. So please use OpenCV with Python. Please check out the [supplemental document](#) for instructions on using SIFT in VLFeat/MATLAB or OpenCV/Python.

The standard use of a SIFT library consists of you just providing an image and the library does its thing: finds interest points at various scales and computes descriptors at each point. We're going to use the library code only to compute the orientation histogram descriptors for the interest points you have already

detected from Problem 1. To do so, you need to provide a scale setting and an orientation for each feature point as well as the gradient magnitude and angle for each pixel. The scale we'll fix to 1.0 (see accompanying SIFT software usage tutorial). The orientation needs to be computed from the gradients:

$$g_{\theta} = \text{atan2}(I_y, I_x)$$

But, you already have the gradient images! So you can create an "angle" image and then for a given feature point at  $\langle u_i, v_i \rangle$  you can get the gradient direction.

a. Write the function to compute the angle. Then for the set of interest points you found above, plot the points for all of transA, transB, simA and simB on the respective images and draw a little line that shows the direction of the gradient. In MATLAB if you want you can use the VLFeat function `vl_plotframe` to draw the feature points locations and the angle. You'll need to figure it out - look at <http://www.vlfeat.org/overview/sift.html> and also the documentation for `vl_plotframe`. In OpenCV you can use the method `drawKeypoints()`.

**Output:**

-Interest points with angles shown on (transA, transB) pair as ps4-2-a-1.png

-Interest points with angles shown on (simA, simB) pair as ps4-2-a-2.png

b. Now we're going to call the SIFT descriptor code. You need to pass in each keypoint location along with its scale and orientation. This process is covered in more detail in the accompanying supplemental document.

Once we have the descriptors, we need to match them. This is what we called *putative* matches in class. Given keypoints in two images, get the best matches. Both VLFeat and OpenCV have functions for computing matches (e.g. for VLFeat it's called `vl_ubcmatch`.) You will call those and then you will make an image that has both the A and B version adjoined and that draws lines from each keypoint in the left to the matched keypoint in the right. We'll call this new image the putative-pair-image.

Now, both of the SIFT packages have functions to draw matches. **\*\*\*But you are not permitted to use them!!!\*\*\*** This way you will have to identify the location of each keypoint and its match, and explicitly draw the line. This tells us you have a good handle on the data structures that you are using.

Write the function to call the appropriate SIFT descriptor extraction function with the necessary input data structures. Do this for all the keypoints in both pairs of images. Then call the matching functions of VLFeat or OpenCV to compute the best matches between the left and right images of each pair. Then create the putative-pair-image for both transA-transB and simA-simB pair. You must write your own drawing function (note you may use OpenCV's `line()` function or MATLAB's

plot() function).

**Output:**

- putative-pair-image for (transA-transB) as ps4-2-b-1.png
- putative-pair-image for (simA-simB) as ps4-2-b-2.png

### 3. RANSAC

We're almost there. You now have keypoints, descriptors and their putative matches. What remains is RANSAC. To do this for the translation case is easy. Using the matched keypoints for transA and transB, randomly select one of the putative matches. This will give you an offset (a translation in X and Y ) between the two images. Find out how many other putative matches agree with this offset (remember, you may have to account for noise, so "agreeing" means within some tolerance). This is the *consensus* set for the selected first match. Find the best such translation - the one with the biggest consensus set.

- Write the code to do the translational case on transA and transB. Draw the lines on the adjoined images of the biggest consensus set.

**Output:**

- biggest consensus set lines drawn on pair (transA -transB) as ps4-3-a-1.png

**Output (Textual Response):**

-What translation vector was used?

-What percentage of your matches was the biggest consensus set?

- For the other image pair we need to compute a similarity transform. Recall that a similarity transform allows translation, rotation and scaling. We can represent this transform with a matrix as:

$$\begin{bmatrix} u' \\ v' \end{bmatrix} = \begin{bmatrix} a & -b & c \\ b & a & d \end{bmatrix} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$$

In other words, there are four unknowns. Each match gives two equations - so you need to pick two matches to solve, which is why similarity is called a two point transform. Clever.

Do the same as above but for the similarity pair simA and simB. Write code to apply RANSAC by randomly picking two matches, solving for the transform, and determining the consensus set. Draw the lines on the adjoined images for the biggest consensus set.

**Output:**

- biggest consensus set lines drawn on pair (simA -simB) as ps4-3-b-1.png

**Output (Textual Response):**

- What is the transform matrix for the best set?
- What percentage of your matches was the biggest consensus set?

**EXTRA CREDIT QUESTIONS (3-c, 3-d, 3-e)**

c. For the second image pair we told you that the transform to compute was a similarity transform. But suppose you didn't know that. You might have guessed that it was an affine transform, expressed as follows:

$$\begin{bmatrix} u' \\ v' \end{bmatrix} = \begin{bmatrix} a & b & c \\ d & e & f \end{bmatrix} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$$

Here there are six unknowns. Again, each match gives two equations - so this time you need to pick three matches to solve, which is why affine is called a three point transform. Still clever?

Try estimating the affine transform between simA and simB. Write code to apply RANSAC by randomly picking three matches, solving for the transform, and determining the consensus set. Draw the lines on the adjoined images for the biggest consensus set.

**Output:**

- biggest consensus set lines drawn on pair (simA -simB) as ps4-3-c-1.png

**Output (Textual Response):**

- What is the transform matrix for the best set?
- What percentage of your matches was the biggest consensus set?

d. Finally, given these transforms, you should be able to warp the second image to the first. We're not going to tell you about how to that here except we did talk about warping (remember backward warping?) earlier.

Create a new version of simA by warping simB "back" to the coordinate system of simA using the 3-b transform you found. Call the new image warpedB. Show the two images (simA and warpedB overlaid by either blending them or by making a pseudo color image where you put simA in the red channel and warpedB in the green channel of a color image (both (simA and warpedB are grayscale images).

**Output:**

-warpedB image as ps4-3-d-1.png  
-the overlay image as ps4-3-d-2.png  
e. Do 3-d again but this time using the affine transform recovered in 3-c.

**Output:**

-warpedB image as ps4-3-e-1.png  
-the overlay image as ps4-3-e-2.png

**Output (Textual Response):**

Comment as to whether using the similarity transform or the affine one gave better results, and why or why not.