

HOMEWORK #4

University of Southern California

Lecturer: C.-C.JayKuo

Issued: 3/4/2019

Due:3/19/2019

Report

1 Problem 1: Texture Analysis

1.1 Texture Classification

1.1.1 Abstract and Motivation

Texture is a description of the spatial arrangement of color or intensities in an image or a selected region of an image. Laws developed a series of filters that could automatically detect the texture patterns. Laws' masks are well established as one of the best methods for texture analysis in image processing and are used in various applications, including medical image analysis. In this section, I use the five 1D laws filters to extract the features of 12 textures images and use kmeans method to classification.

1.1.2 Approach and Procedures

1. Preprocessing

First of all, extend the image boundary with 2 pixels width for each edge by reflection. Then, normalize the input image pixel value by,

$$I(i, j) = \frac{I(i, j)}{255} \quad (1)$$

Use the result image as the input for the following steps.

2. Set up the Laws filter

The 2D Laws' filter is caculated by using 1D filter. The five 1D filters are as follows.

$$\begin{aligned} L5 &= [1, 4, 6, 4, 1] \\ E5 &= [-1, -2, 0, 2, 1] \\ S5 &= [-1, 0, 2, 0, -1] \\ W5 &= [-1, 2, 0, -2, 1] \\ R5 &= [1, -4, 6, -4, 1] \end{aligned} \quad (2)$$

This filter detect level, edge, spot, wave and ripple of the input image. The 2D Laws' filter is the outer product of the 1D laws filters. Thus, we can finally get 25 2D Laws' filters. For one Laws filter, it detect the different property in different direction. Take the E5S5 as an example, the E5S5 is

$$E5S5 = \begin{bmatrix} 1 & 0 & -2 & 0 & 1 \\ 2 & 0 & -4 & 0 & 2 \\ 0 & 0 & 0 & 0 & 0 \\ -2 & 0 & 4 & 0 & -2 \\ -1 & 0 & 2 & 0 & -1 \end{bmatrix} \quad (3)$$

In the horizontal direction, it detects the edge characteristic of the image, while in the vertical direction, it detects the spot characteristic.

3. Convolution and Calculate the Energy

Use the 2D filter obtained above, convolute the input image and get a filtered image. Calculate the square sum as the representation of the energy using the following equation,

$$Square_sum = \sqrt{\frac{1}{MN} \sum_{i=1}^{i=M} \sum_{j=1}^{j=N} I(i, j)^2} \quad (4)$$

Where, M is the height and N is the width of the image. The $Square_sum$ is represent the texture property obtained by each filter. After filtered with 25 2-D laws filters, we can finally get a 1 by 25 dimension vector for each images. Since we have 12 texture images, we finally get a 12×25 two-dimension array.

4. Principal Component Analysis

Suppose the input matrix is a $N \times D$ matrix, where each row represent a training point while each column represent the feature of the data.

First, normalize the feature of each data into zero mean and unit variance, i.e. subtract the mean from each column and divided its variance.

$$X(:, i) = \frac{X(:, i) - \mu}{var} \quad (5)$$

Second, calculate the covariance matrix \mathbf{C} .

$$cov(\mathbf{C}) = \mathbf{X}^T \mathbf{X} \quad (6)$$

Third, do the SVD decomposition of \mathbf{C} .

$$\mathbf{C} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T = \sum_{i=1}^{i=D} u_i \sigma_i v_i \quad (7)$$

Where $\sigma_1, \sigma_2 \dots \sigma_D$ is the singular value sorted in decreasing order. Pick the first three largest singular value and the corresponding left and right singular vectors to form the new low dimension dataset.

Last, recover the low-dimension data by multiply the first three columns of the left-singular-vector $\mathbf{U}(:, 1 : 3)$.

$$\hat{\mathbf{X}} = \mathbf{X} \times \mathbf{U}(:, 1 : 3) \quad (8)$$

5. K-means Clustering

For a dataset $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$, the K-means distortion objective function is,

$$J(\{\mu_k\}, \{r_{ik}\}) = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K r_{ik} \| \mu_k - \mathbf{x}_i \|_2^2 \quad (9)$$

Where μ_1, \dots, μ_N are the centroids of K clusters and $r_{ik} \in \{0, 1\}$ represent whether example i belongs to cluster k . Our goal is to minimize J by update the centroid and clusters over iteration.

First, fixing the centroid and cluster the dataset. The cluster clasification is measured by Euclidean distance.

$$\hat{r}_{ik} = \begin{cases} 1 & k = \operatorname{argmin}_{k'} \|\mu_{k'} - \mathbf{x}_i\|_2^2 \\ 0 & \text{else} \end{cases} \quad (10)$$

Second, fixing the cluster assignment and find the new centroid by minimize J .

$$\hat{\mu}_k = \frac{\sum_{i=1}^N r_{ik} \mathbf{x}_i}{\sum_{i=1}^N r_{ik}} \quad (11)$$

By repeating this two steps, k-means finally cluster the input data into k cluster until the maximum iteration is reached.

1.1.3 Experimental Result

The feature obtained by laws filters are shown in Table 1.

Table 1: Laws' filter Results

	L5L5	L5E5	L5S5	L5W5	L5R5	E5L5	E5E5	E5S5	E5W5	E5R5	S5L5	S5E5	S5S5	S5W5	S5R5	W5L5	W5E5	W5S5	W5W5	W5R5	R5L5	R5E5	R5S5	R5W5	R5R5	
Texture 1	23.234247	4.6825110	2.1076640	1.8063850	2.8931140	4.6563900	1.14569	0.615261	0.543388	0.911480	2.236050	0.629166	0.366846	0.348021	0.596002	1.8570570	0.571519	0.353790	0.344430	0.624239	3.0070380	0.992611	0.656907	0.660362	1.2684020	
Texture 2	15.151940	3.3037110	2.0008330	1.9259510	3.3339640	3.0939180	1.0554210	0.705943	0.718552	0.3045450	1.9633920	0.664501	0.458605	0.483221	0.948022	1.8790030	0.605036	0.465297	0.504723	1.027850	4.22519	1.2137740	0.912609	1.0329280	2.1175569	
Texture 3	16.203750	3.3037110	2.0008330	1.9259510	3.3339640	3.0939180	1.0554210	0.705943	0.718552	0.3045450	1.9633920	0.664501	0.458605	0.483221	0.948022	1.8790030	0.605036	0.465297	0.504723	1.027850	4.22519	1.2137740	0.912609	1.0329280	2.1175569	
Texture 4	23.468594	6.7595880	3.3780370	6.6511910	6.1144160	4.5740730	0.4889290	0.6393800	0.540918	0.53297	2.6184800	0.774882	0.488874	0.479701	0.875657	2.416080	0.724981	0.458605	0.483221	0.948022	1.8790030	0.605036	0.465297	0.504723	1.027850	
Texture 5	42.057818	9.6580230	4.8892210	3.9940350	6.3130140	2.4742920	1.3236740	1.1989280	0.737567	1.3297530	0.701888	4.7171120	1.3597530	0.701888	0.737567	1.2900120	4.00038	1.2359560	0.764760	0.7422110	1.3488630	7.1464640	2.3008030	1.4254330	1.4456380	2.7380820
Texture 6	37.257572	10.672662	6.2258550	5.6340540	9.3717850	7.3090420	2.0258610	1.3637470	1.3177120	2.4033310	3.8568120	1.2125390	0.763727	0.750567	1.4043120	3.5646210	1.1379410	0.725143	0.722346	1.3721710	6.26047	2.0481660	1.3295620	1.3556710	1.2558960	
Texture 7	22.204380	5.1025030	2.5982420	2.1287460	3.3601380	4.8471780	1.9332860	0.717721	0.637370	1.0068270	2.5244110	0.712157	0.419269	0.390548	0.683647	2.1683510	0.653331	0.403365	0.393021	0.713837	4.1362860	1.1646540	0.756777	0.767413	1.44918	
Texture 8	29.347879	6.2740220	3.7886010	3.6377180	6.3145750	5.7985960	2.0652240	1.3701370	1.3871990	2.5091330	3.4533610	1.3898360	0.905245	0.949446	1.888910	3.2782250	1.2690810	0.929327	1.00544	0.2069070	6.2573210	2.3643810	1.8114760	2.0633110	4.2256870	
Texture 9	43.238528	9.1141560	6.7507420	7.3943790	13.912112	3.7394160	0.563597	0.582703	1.1947090	1.9297210	0.234479	0.395754	0.355980	0.789401	1.0991890	1.9029960	0.536469	0.353914	0.375567	0.789401	4.2254170	1.0991890	0.736088	0.796160	1.6858169	
Texture 10	23.015021	4.8893070	3.2974780	3.425269	6.6216140	2.3702640	0.976126	0.759962	0.889601	1.9975260	1.0646110	0.428021	0.343537	0.423066	1.0391950	0.890206	0.327091	0.256623	0.319995	0.834625	1.1443550	0.505643	0.387378	0.4711	1.2521780	
Texture 11	16.264758	4.8347470	3.1738660	3.0221600	4.9084640	1.3776630	0.332087	0.217935	0.234377	0.495369	0.179485	0.188805	0.125442	0.137233	0.304436	0.780822	0.187067	0.126896	0.1465	0.310427	1.5262950	0.379699	0.269653	0.293954	0.662665	
Texture 12	26.865731	6.6678070	3.6361530	3.2245420	5.4255880	5.4547440	1.58110	0.773449	0.740840	1.3593010	3.3570780	0.697736	0.423219	0.415515	0.790853	2.102330	0.640086	0.395435	0.395456	0.766898	3.35808	1.1281490	0.721796	0.744824	1.4777550	

The plot of PCA is shown in Figure 1.

The kmeans result are as follows.

Groud truth: Bark: 4, 6, 12. Straw: 2, 8, 10. Brick: 3, 9, 11. Bubble: 1, 5, 7.

Use 25 dimension feature space to do the kmeans: 1, 2, 3, 1, 4, 4, 1, 1, 4, 1, 3, 1.

Use 3 dimension feature space to do the kmeans: 2, 2, 3, 1, 4, 4, 2, 4, 1, 2, 3, 1.

1.1.4 Discussion

Compute the variance of each feature over 12 images, we can conclude that L5L5 maintains the largest discriminant power, that could distinguish each images. That is because the L5L5 filter collects the low-dimension information of the input image, which is dominant in the image. Beside, the Numerical value is the highest among the features. Thus, a small fluctuation of L5L5 will result in a large variance in the feature space. On the other hand, S5W5 and S5S5 maintains the smallest discriminant power, because they extract the high frequency information which is of small quantity in the images.

Table 2: Variance of the 25-D features

	L5L5	L5E5	L5S5	L5W5	L5R5	E5L5	E5E5	E5S5	E5W5	E5R5	S5L5	S5E5	S5S5	S5W5	S5R5	W5L5	W5E5	W5S5	W5W5	W5R5	R5L5	R5E5	R5S5	R5W5	R5R5
var	92.6800	5.3703	2.1778	2.4890	9.1832	5.1849	0.4499	0.1569	0.1422	0.4579	1.3863	0.1519	0.0599	0.0583	0.2046	1.0529	0.1383	0.0609	0.0637	0.2342	3.3524	0.4508	0.2239	0.2642	1.0159

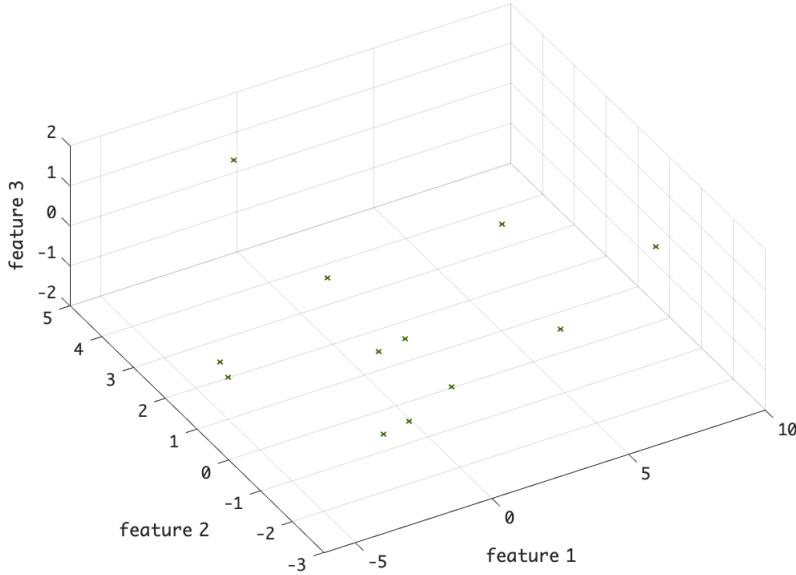


Figure 1: PCA Result

The kmeans method do not work very well in texture classification. With 25D feature space, it only cluster “brick” successfully and some “bubbles”. It makes some mistakes that it cluster “straws” into “bubbles”. That might because the feature representation of this two textures are not strong enough. With 3D feature, the result improves a little. Althought it still cannot seperate the “straws” and “bubbles”, but it cluster the “Bark” and “brick” correctly. That might because the “brick” has strong vertical edge and horizontal edge information that distinguish among the textures pattern while the “Bark” also maintains strong vertical edge information.

1.2 Texture Segmentation

1.2.1 Abstract and Motivation

Base on the previous section, we use kmeans method to classify different feature base on the feature property. In this section, we use the same method – Apply 25 Laws’ filter to segment the image contains various feature.

1.2.2 Approach and Procedures

1. Subtract the Local Mean

First, subtract the local mean from each pixel within a 5 by 5 window to eliminate the brightness in the same texture.

$$I(\hat{i}, j) = I(i, j) - \frac{1}{N \times N} \sum \sum I(k, l) \quad (12)$$

2. Feature extraction

As it is describe in 1.1.2, apply 25 2-D Laws filters to extract the 25 feature of the entire image pixels. After this step, I got a $Row \times Column \times 25$ three-dimension array.

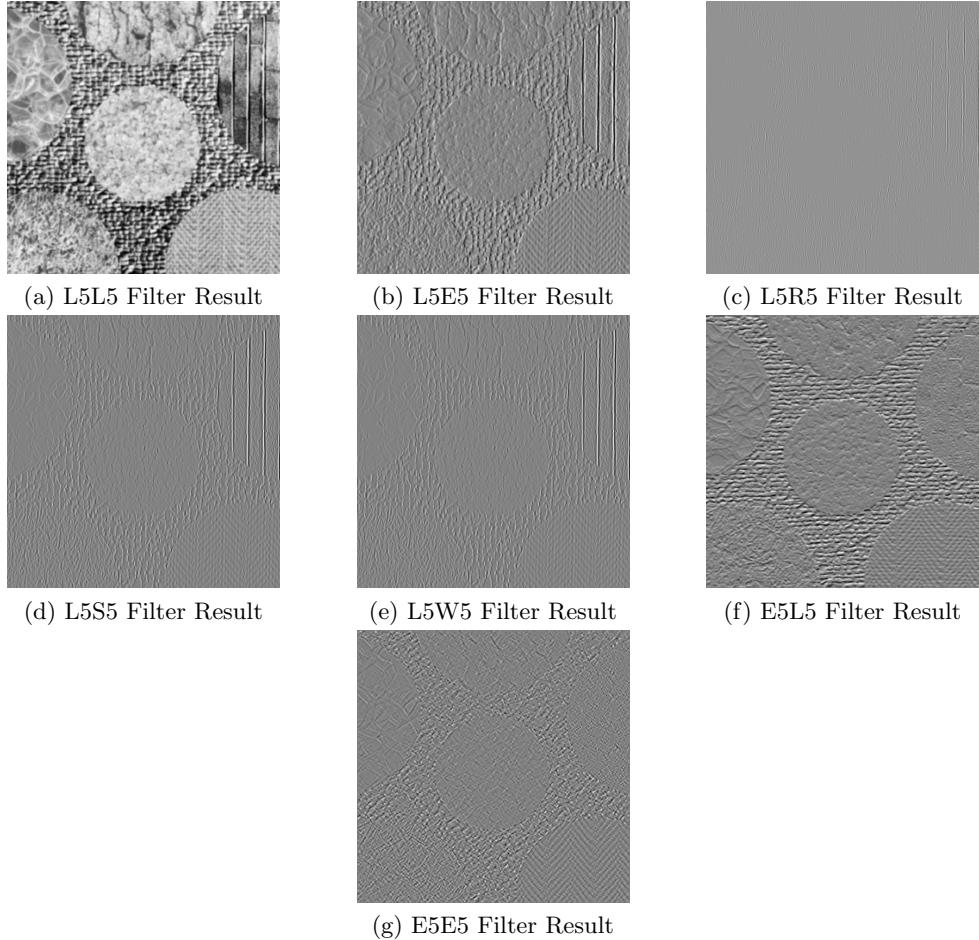


Figure 2: Some Filter Result

3. Compute Energy Map

Use the Window to produce 25-d energy map.

$$Square_sum = \sqrt{\frac{1}{M^2} \sum_{i=1}^{i=M} \sum_{j=1}^{j=M} I(i,j)^2} \quad (13)$$

Where, M is the window size. Here, I choose $M = 7, 15, 17$. Use this value to represent the energy of the pixel located at the center of the window.

4. K-means cluster

Apply k-means cluster with 25-D feature space. Set the number of cluster equals to 7. Output the result of the segmented regions by using 7 different grayscale intensity.

1.2.3 Experimental Result

The segmentation result is shown in Figure 3.



Figure 3: Segmentation Result, Window size = 15

1.2.4 Discussion

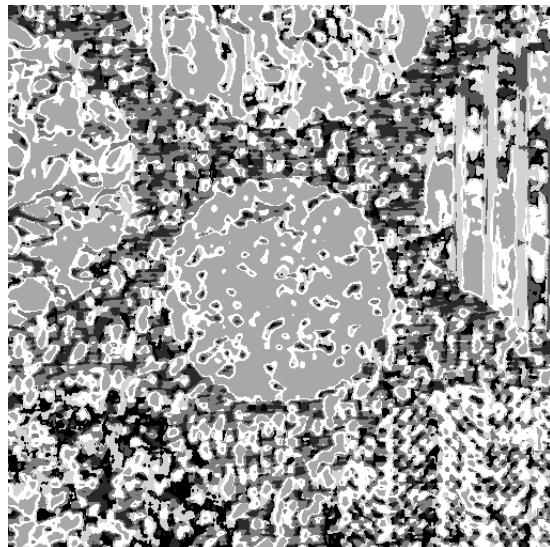


Figure 4: Segmentation Result ,Window size = 7



Figure 5: Segmentation Result ,Window size = 17

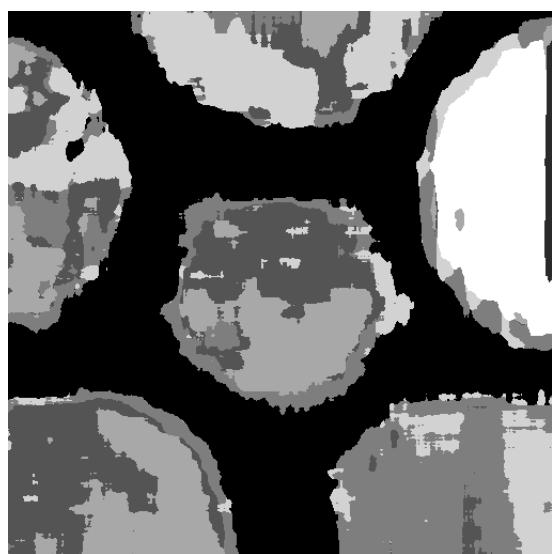


Figure 6: Segmentation Result ,Window size = 41

1.2.5 Discussion

Figure 4 to 6 shows the segmentation result is not quite good. From the figure, we can observe that the larger window size is, the less isolate holes remain in a particular regions. When window size = 41, it shows the best result. Because it separates the most areas. There still remains some holes in the image and the texture are not correctly classified. That might because the 25-D feature are not strong enough to represent the pixel information. Thus, PCA will not help much in this case.

1.3 Advanced Texture Segmentation Techniques

1.3.1 Abstract and Motivation

The result of 1(b) is not so clean and in this section I will come up with several idea to enhance the results.

1.3.2 Approach and Procedures

1. Adopt PCA

Adapt a PCA method, reduce the 25-D feature space to low dimension feature space. Use the low dimension feature space to do the k-means and output the segmentation result. I try to reduce the dimension to 15, 7 and 3 and observe the result.

2. Post-processing

Adapt median filter to merge the small hole and make the segmentation result more cleaning. The main idea of the median filter is to run through the signal entry by entry, replacing each entry with the median of neighbourhood entries. The pattern of neighbors is called the “window”, which slides, entry by entry, over the entire signal. For 2D images, the window size normally has an odd number of entries, like 3, 5 or 7. Here I choose the window size equal to 5. I filter the image for 30 times and the result are shown in the next section.

3. Enhance the Boundary

Looking at boundary, the pattern of the boundary pixels are unique. First, itself is not the background color(which in this case is 0), then, it's 4-connected neighbour must contain a background color. Collect the boundary pixels of a texture pattern. Use the majority vote method to decide the boundary pixels' class.

1.3.3 Experimental Result

1.3.4 Discussion

As I mentioned above, the PCA method does not help much in improving the segmentation result. However, when taking a close look at the detail, **PCA does merge some small holes in to a continuous area**. But it still enlarge the boundary around the texture pattern. Although the segmentation result is not so good, but the median filter actually work well in merge the small hole in a specific area. From the Figure 10, we can see that the boundary is more clearly. And the small hole in the texture pattern are merged. The k-means method works well in classify the texture located at the right-down area and the middle area.

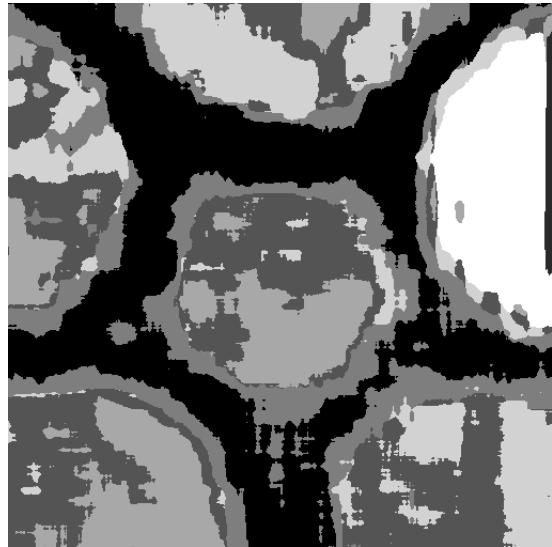


Figure 7: PCA to 3-D Result, Window size = 41

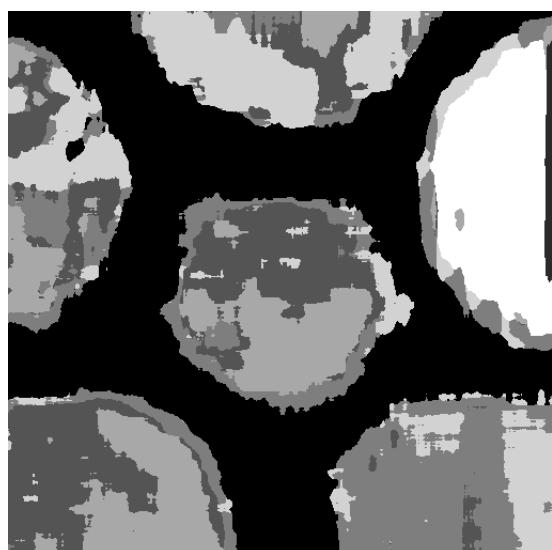


Figure 8: PCA to 7-D Result, Window size = 41

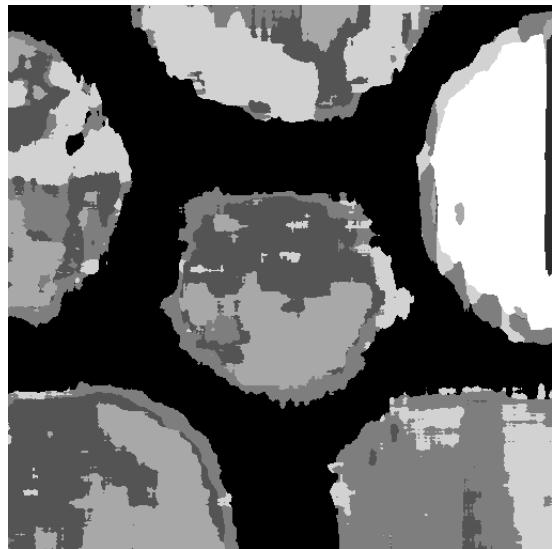


Figure 9: PCA to 15-D Result, Window size = 41

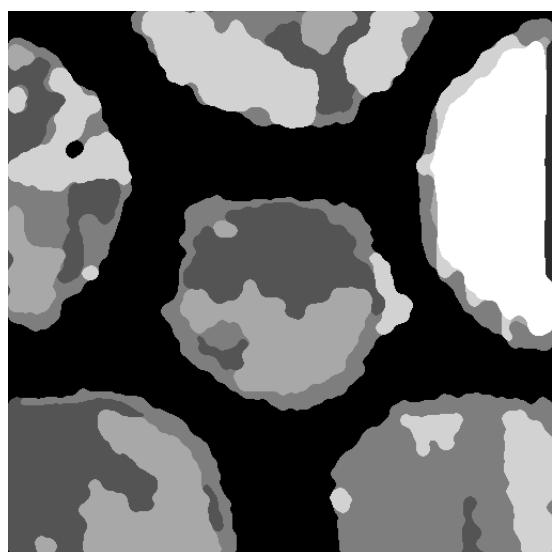


Figure 10: Median Filter Result, Window size = 41

2 Problem 2: Image Feature Extractor

2.1 SIFT

2.1.1 Abstract and Motivation

Scale-invariant feature transform (SIFT) is a computer vision algorithm used to detect and describe local features in image. It looks for the points that invariant in **translation, scaling and rotation** in spatial scale. This algorithm was published by David Lowe in 1999, and was summarized in 2004. In this section, I will introduce the basic procedure of the SIFT algorithm.

2.1.2 Approach and Procedures

The fundamental of the SIFT is as follow. The scaling-invariance and position-invariance is achieved by selecting key location at the maxima and minima of the DoG applied in scale space. This key location are treated as the potential key points. First, the input image is continuously filtered and downsampled by stage Gaussian kernel functions with different σ and building an image pyramid. After that, we need to remove the points that are not rotation-invariant. It mainly includes the following key steps.

1. Feature extraction

In order to make the features have **position and scale invariance**, the detection of feature points is done in image scale space. The input image is convolved with Gaussian function through several level.

$$\begin{aligned} G(\sigma_1) &= (2\pi\sigma_1^2)^{-1}e^{-(x^2+y^2)/2\sigma_1^2} \\ G(\sigma_2) &= (2\pi\sigma_2^2)^{-1}e^{-(x^2+y^2)/2\sigma_2^2} \\ G(\sigma_3) &= (2\pi\sigma_3^2)^{-1}e^{-(x^2+y^2)/2\sigma_3^2} \\ G(\sigma_4) &= (2\pi\sigma_4^2)^{-1}e^{-(x^2+y^2)/2\sigma_4^2} \\ G(\sigma_5) &= (2\pi\sigma_5^2)^{-1}e^{-(x^2+y^2)/2\sigma_5^2} \end{aligned} \quad (14)$$

Where $\sigma_1 = \sqrt{2}$, $\sigma_2 = 2$, $\sigma_3 = 2\sqrt{2}$, $\sigma_4 = 4$, $\sigma_5 = 4\sqrt{2}$.

Sample the image with different size. For each specific size, the image are filtered by Gaussian filters with different σ . In order to detect stable and unique key points **more efficiently**, then subtract the adjacent gaussian smooth output image to approximate the Laplacian of Gaussian (LoG) with in the same octave. **The computation of DoG cost less than the Computation of LoG.** Then, for each pixels, we choose the 26 neighbourhood pixels to compare with. If it is the minima or maxima of the neighbourhood, we regard it as a potential key point.

2. Feature description

After get these potential key points, the gradient direction distribution characteristic of the key point's neighborhood pixels is used to assign a direction to each image descriptors, so that the key point has **rotation invariance**.

$$M_{ij} = \sqrt{(A_{i,j} - A_{i+1,j})^2 + (A_{i,j} - A_{i,j+1})^2} \quad (15)$$

$$R_{i,j} = \tan^{-1} (A_{i,j} + A_{i+1,j}, A_{i,j+1} + A_{i,j}) \quad (16)$$

The pixel differences are efficient to compute and provide sufficient accuracy due to the substantial level of previous smoothing.

3. Removing

By removing the potential points that not possess the rotation-invariance, the key point achieves its robustness of rotation-invariant.

(a) Detect “Low contrast point”

Use Taylor expansion method to compute the contrast. Compute the Hessian matrix \mathbf{H} and its trace and determination.

$$\mathbf{H} = \begin{bmatrix} D_{xx} & D_{xy} \\ D_{xy} & D_{yy} \end{bmatrix} \quad (17)$$

$$tr(\mathbf{H}) = 2D_{xy} \quad (18)$$

$$det(\mathbf{H}) = D_{xy}^2 - D_{xx}D_{yy} \quad (19)$$

if

$$\frac{tr(\mathbf{H})^2}{det(\mathbf{H})} > \frac{(r+1)^2}{r} \quad (20)$$

remove it.

(b) Remove points in non-salient region

In order to make it as stable as possible, the orientation is determined by the peak in the histogram of local imgae gradient orientation. The orientation histogram has 36 bins covering 360 degree range of orientation, ans select the peak value to represent the orientation of the key point.

4. Illuminance Change

SIFT algorithm guarantee its robustness to illuminance change by thresholding the gradient magnitudes with maximum possible gradient value times 0.1.

5. Output descriptor

The orientation planes is formed by collecting the pixels that fall in a circle of radius 8 pixels around a key location. The orientation is defined by subtracting the key-points’s orientation. For each key-points, use 8 orientation planes to sample over a 4 by 4 grid of locations, with a sample spacing 4 times that of the pixel spacing used for gradient detection. In order to sample the image at larger scale, the same process is repeated for a second level. But in this time, we use a 2 by 2 region. Thus, we got $8 \times 4 \times 4 + 8 \times 2 \times 2 = 160$ elements. SIFT’s output vector size in its original paper contains 160 elements.

2.2 Image Matching

2.2.1 Abstract and Motivation

Image retrieval is an important task in object recogonization area. In this section, I will implement the SIFT algorithm on OpenCV toolbox to extract the feature vector and do the Nearest Neightbour Search to do the image retrieval.

2.2.2 Approach and Procedures

1. Find the Key-point & Extract SIFT Feature

Input the image *river1.raw* and *river2.raw* and extract the SIFT feature of the two images. The number of key points found by OpenCV_SIFT is shown in Figure 11. The feature extract by OpenCV_SIFT contains 128 dimensions.

```
number of key point extract from river1.raw: 8618
number of key point extract from river2.raw: 9082
```

Figure 11: number of key point found by OpenCV_SIFT

2. Find the Largest Scale Key-point

The “largest scale key point” is defined by the key point’s corresponding descriptor with the largest L_2 norm. Thus, by traversing the entire SIFT descriptors of the *river1.raw* and caculating the L_2 norm, I found the largest scale key point \mathbf{kp}_1 . The parameters of \mathbf{kp}_1 is shown in Figure 12. The position of \mathbf{kp}_1 is marked in Figure 13.

```
The index of the largest scale key point in river1.raw: (60.05622100830078, 476.3304443359375)
The orientation of the largest scale key point in river1.raw: 206.17921447753906
The diameter of the largest scale key point neighbourhood in river1.raw: 6.630932807922363
```

Figure 12: The parameters of \mathbf{kp}_1



Figure 13: The position of \mathbf{kp}_1

3. Find the Closest Neighboring Key-point Use K-NN method and compute the L_2 norm of the distance between \mathbf{kp}_1 and the key-points found in *river2.raw*. The L_2 norm is defined by

$$\|\mathbf{kp}_1 - \mathbf{kp}'\|_2 = \sqrt{(x_1 - x'_1)^2 + (x_2 - x'_2)^2 + \dots + (x_D - x'_D)^2} \quad (21)$$

Where \mathbf{kp}' is the key-point of *river2.raw*. The $x_1, x_2 \dots x_D$ are $D = 128$ features of the descriptor, i.e. the 128 entries in a row of the descriptor.

The key-point with the shortest $L2$ norm distance is regarded as the closest neighboring key-point of \mathbf{kp}_1 . I try several number of nearest neighbour and the result are shown in the next section.

2.2.3 Experimental Result

The matching result of *river1.raw* and *river2.raw* is shown in Figure 14, where I choose 500 key-points in each image. The Nearest Neighbour in *river1.raw* and its orientation is circled in Figure 15. The parameter of the key-point is shown in Figure 16

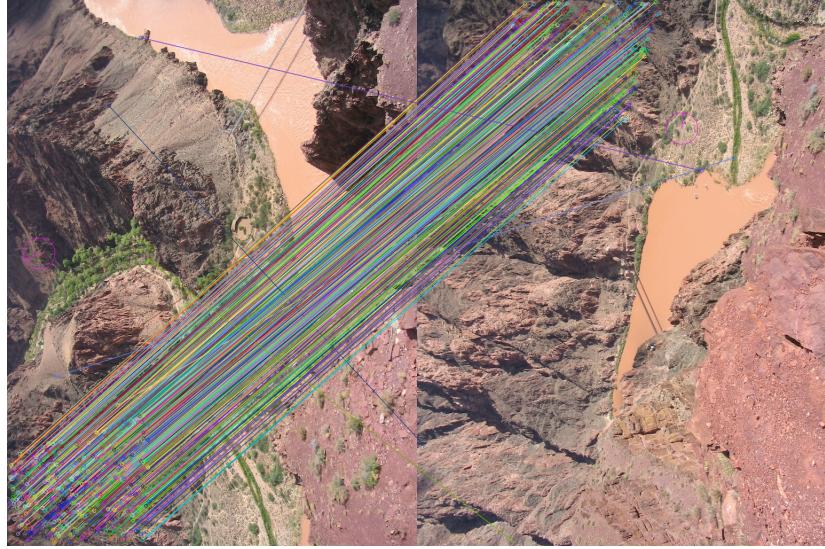


Figure 14: 500 key points Matching Result

The 2-Nearest Neighbour in *river1.raw* and its orientation is circled in Figure 17. The parameter of the key-point is shown in Figure 18.

The 3-Nearest Neighbour in *river1.raw* and its orientation is circled in Figure 19. The parameter of the key-point is shown in Figure 20

2.2.4 Discussion

The dominant orientation of the nearest point in *river2.raw* is 129.87, which is different from the \mathbf{kp}' . Though it is true that in descriptor feature the dominant orientation is provided, this is not used to actually compare the similarity of the keypoints. The keypoints are primarily compared using the distance between the descriptor, which regardless of this dominant orientation are inherently orientation invariant. The dominant orientation are used in some retrieval/classification tasks as extra information used to improve matching scores.



Figure 15: 1-NN result

```
nearest point in river2: (28.139734268188477, 466.1160583496094)
the orientation of the nearest point in river2: 129.87998962402344
The diameter of of the largest scale key point neighbourhood in river1.raw: 1.9627506732940674
```

Figure 16: Parameter of the 1-NN Result



Figure 17: 2-NN result

```
the orientation of the nearest point in river2: 129.87998962402344
The diameter of of the largest scale key point neighbourhood in river1.raw: 1.9627506732940674
nearest point in river2: (496.109375, 242.12303161621094)
the orientation of the nearest point in river2: 65.8524169921875
The diameter of of the largest scale key point neighbourhood in river1.raw: 9.828125953674316
```

Figure 18: Parameter of the 2-NN Result



Figure 19: 3-NN result

```
the orientation of the nearest point in river2: 129.87998962402344
The diameter of the largest scale key point neighbourhood in river1.raw: 1.9627506732940674
nearest point in river2: (496.109375, 242.12303161621094)
the orientation of the nearest point in river2: 65.8524169921875
The diameter of the largest scale key point neighbourhood in river1.raw: 9.828125953674316
nearest point in river2: (411.07147216796875, 162.6121826171875)
the orientation of the nearest point in river2: 315.71728515625
The diameter of the largest scale key point neighbourhood in river1.raw: 1.9489965438842773
```

Figure 20: Parameter of the 3-NN Result

2.3 Bag of Words

2.3.1 Abstract and Motivation

A Bag-of-Words model that encodes in graphs the local structures of a digital object. In this section, I use the samples of zeros and ones from MNIST to form a codebook with bin size of 2. Use the SIFT algorithm to extract the feature from the image and do k-means cluster. Then, do the classification for the *eight.raw* and show the histogram.

2.3.2 Approach and Procedures

1. Extract the SIFT feature

Extract the SIFT feature for each image *one_1.raw* ... *zero_5.raw*. Collect the descriptor for all the image and form a training data set. The OpenCV_SIFT find 54 key point in total. Thus, the training dataset is a 54×128 2-D array.

```
number of key point in 0: 8
number of key point in 0: 2
number of key point in 0: 12
number of key point in 0: 12
number of key point in 0: 7
number of key point in 1: 2
number of key point in 1: 3
number of key point in 1: 3
number of key point in 1: 0
number of key point in 1: 5
```

Figure 21: number of key-points

2. K-menas Cluster

Train the training set with k-means cluster method and divide the training point into 2 clusters. The zeros' key points are shown in Figure 22 while the ones' key-points are shown in Figure 23. The blue dot represent the cluster of object, while the magenta dot represent the cluster of background.

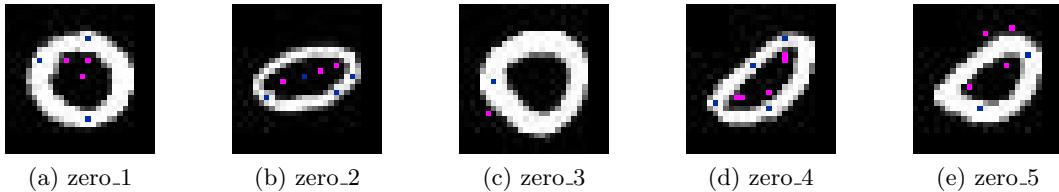


Figure 22: Zeros key-points

Finally we got 2 centroid and 2 cluster for futher prediction.

3. Get the Histogram

With the trained codebook, I can use it to get the histogram of the *eight.raw*. Extract the key-points of *eight.raw* and use the discriptor as the prediction dataset. Use the trained k-means cluster above to predict the cluster of the key-points. That is, caculated the $L2$ norm

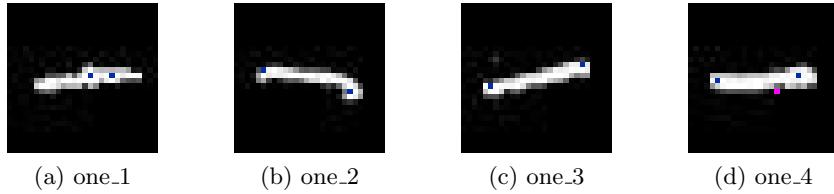


Figure 23: Ones key-points

distance between the key-points and two centroids. Label the key-points with the shortest $L2$ norm distance of the centroid.

4. Caculate the Distance

Regarding the histogram as a 2-dimension vector, caculate the distance between *eight.raw* and the codebook image. First, normalize the vector so that it has unit length. Then caculate the $L2$ norm of the distance and take an average. Use this distance information to classify the target image.

2.3.3 Experimental Result

OpenCV_SIFT extract 12 key points of the *eight.raw*. The cluster result are shown in Figure 24. The histogram of *eight.raw* is shown in 25. The histogram of Zeros is shown in 27. The histogram of Ones is shown in ??.

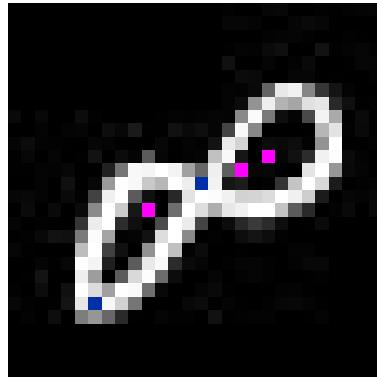


Figure 24: image of eight key-point

2.3.4 Discussion

From the Table 3, we can result that the *eight.raw* is closer to image zeros than ones. If there is only two digits in the bag, then *eight.raw* will be classified as Zeros.

The OpenCV_SIFT detect 12 points in *eight.raw*, however, only 5 of them are marked in the Figure 24. The reason is that some key-points **share the same coordinate but they are different in angle**.

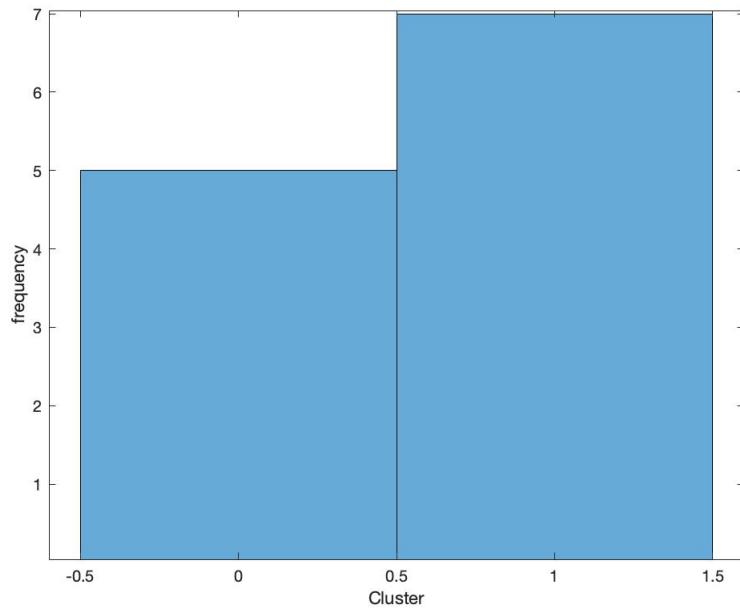


Figure 25: Histogram of *eight.raw*

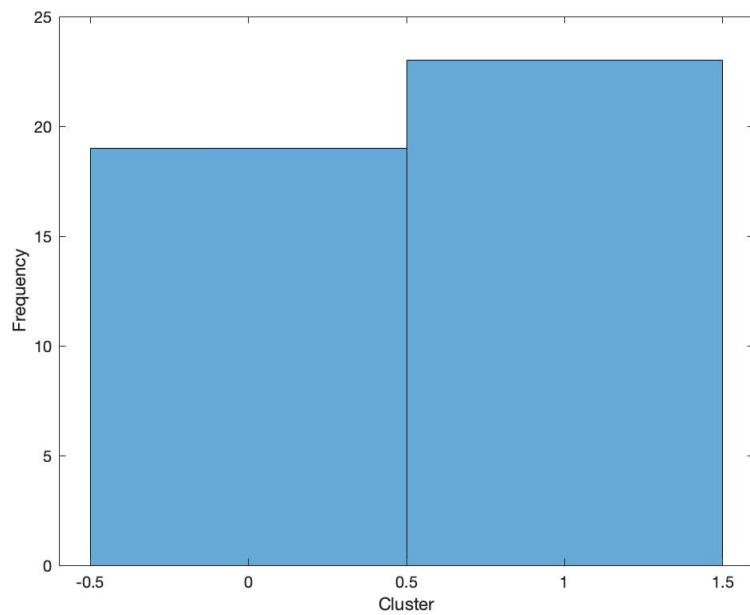


Figure 26: Histogram of Zeros

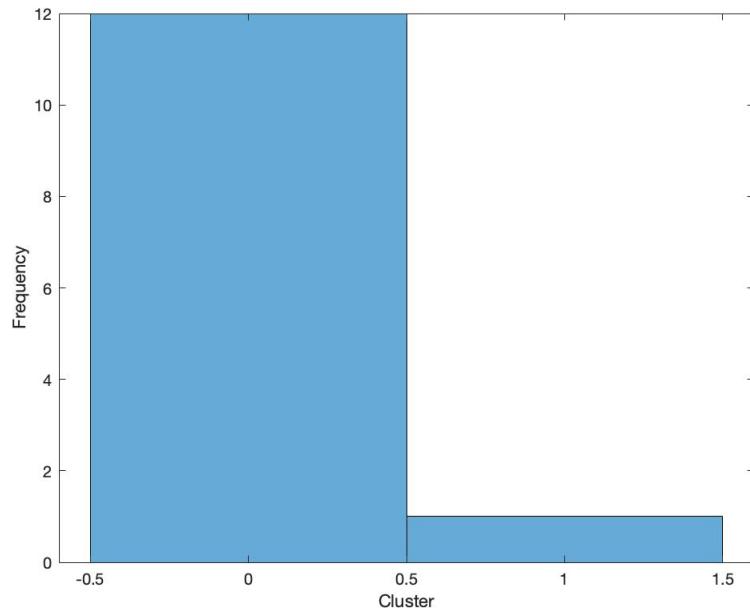


Figure 27: Histogram of Ones

Table 3: Distance Between *eight.raw* and bag of words

codebook	distance
zero_1	0
zero_2	0.4073
zero_3	0.1650
zero_4	0.1650
zero_5	0.2392
one_1	0.9152
one_2	0.9152
one_3	0.9152
one_4	0.6910

```
The index of key-point is: (6.036602020263672, 22.00472068786621)
The orientation of key-point is: 356.158203125
The index of key-point is: (10.599776268005371, 15.82100772857666)
The orientation of key-point is: 62.253753662109375
The index of key-point is: (10.599776268005371, 15.82100772857666)
The orientation of key-point is: 239.60244750976562
The index of key-point is: (10.69171142578125, 15.001066207885742)
The orientation of key-point is: 28.941009521484375
The index of key-point is: (10.69171142578125, 15.001066207885742)
The orientation of key-point is: 220.13241577148438
The index of key-point is: (14.336299896240234, 13.980986595153809)
The orientation of key-point is: 83.46072387695312
The index of key-point is: (14.336299896240234, 13.980986595153809)
The orientation of key-point is: 255.73301696777344
The index of key-point is: (17.977048873901367, 12.411479949951172)
The orientation of key-point is: 87.02120971679688
The index of key-point is: (17.977048873901367, 12.411479949951172)
The orientation of key-point is: 213.2467041015625
The index of key-point is: (19.008071899414062, 11.843730926513672)
The orientation of key-point is: 44.950286865234375
The index of key-point is: (19.008071899414062, 11.843730926513672)
The orientation of key-point is: 80.05438232421875
The index of key-point is: (19.008071899414062, 11.843730926513672)
The orientation of key-point is: 218.8138427734375
```

Figure 28: Parameters of the key-point in *eight.raw*