

CSCI 5521: Introduction to Machine Learning

Homework 3

Saksham Goel | goelx029 | 5138568

1. **Apply PCA to the digits data set to reduce to dimensions needed to capture 90% of the variance.**

I applied the PCA algorithm over the whole dataset. By whole dataset I mean all the 3000 samples containing all three digits (0, 8, 9) and both types of sample (training and test). After applying PCA on that I found out that 73 dimensional space is enough to capture 90% variance.

2. **Write your own K-means algorithm and apply it to the Digits data set, after reducing the dimension using the PCA in the previous step. Use $k = 6$ clusters and initial centers equal to the elements # 1, 1000, 1001, 2000, 2001, 3000 in the original data set. Print out a confusion matrix showing how many 0's, 8's, 9's there are in each cluster. If there are $k = 6$ clusters, this matrix should be 6×3 . If each cluster were assigned a class based on the majority label among members of the cluster, what would be error rate be?**

For this homework, I implemented the K-Means algorithm. For the K-Means algorithm I have two steps which are called the E and M steps. The E step is often called the assignment step. In the assignment step I calculate the distance of each point with each cluster center and assign the cluster to the particular point for which the distance is least. In my matlab implementation I use a vector of assignments where I save the cluster assigned to each point. After the E step comes the M step in which we calculate new cluster center. For this I use the mean values for each sample point and assign that mean value to the cluster center. I repeat the E & M steps again and again until convergence. Here convergence means when the cluster center do not change after a cycle of the E and M step. After running the K-Means algorithm on the digits samples after they have been projected to the 73 dimensional space. I get the following confusion matrix after the K-Means algorithm with the given initialization:

Cluster # \ Digit	0	8	9
1	18	439	14
2	54	489	14
3	5	20	427
4	7	36	528
5	484	6	4
6	432	10	13

To find the Error rate, we can just find the sum of the smallest two entries in each row of the table

above and divide it by the total number of samples. For this part the error rate we got was:

Error Rate: $0.0670 = 6.7\%$

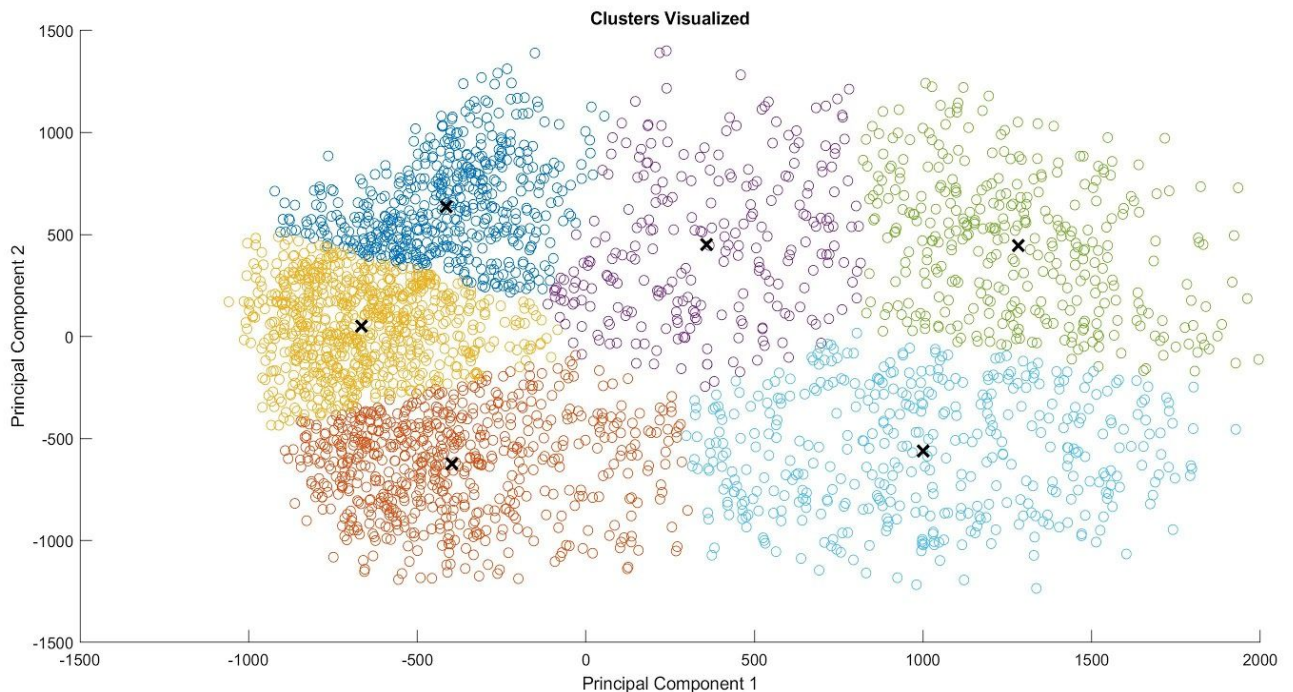
3. **Repeat the above, but start by using only 2 principal components, followed by $k = 6$ clusters. Initialize K-means using the same 6 data samples (projected onto the first two principal components).**

I applied the K-Means algorithm again on the same dataset after the samples were projected on the 2 Dimensional space. The confusion matrix that I got after running K-Means with the given initialization over the dataset after reducing it to a 2 dimensional space is as follows:

Cluster # \ Digit	0	8	9
1	24	368	145
2	37	259	363
3	9	288	455
4	173	71	16
5	355	0	0
6	402	14	21

To find the Error rate, we can just find the sum of the smallest two entries in each row of the table above and divide it by the total number of samples. For this part the error rate we got was:

Error Rate: $0.2947 = 29.47\%$



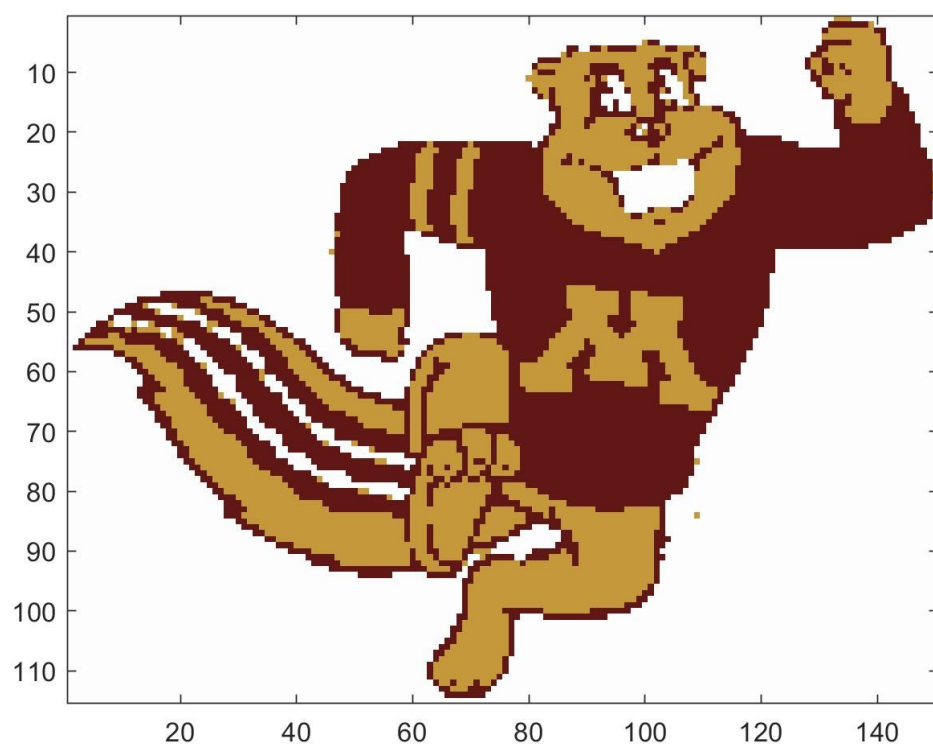
Above is a figure which plots all the points colored based on the clusters they are assigned to. As could be seen that K-Means algorithm did find 6 clusters.

4. **Apply the k-means algorithm to the colored pixel values in image goldy.ppm and stadium.ppm. The data in this case are the RGB pixel values: points in \mathbb{R}^3 . Try $k = 3, 4, 7$. Replace each pixel RGB contents with its corresponding cluster centroid, and re-form the image using the newly substituted pixel values. Redraw the resulting pictures using the modified pixel values. In Matlab, you can read in the picture using the `imread` function, and display it with `imagesc`. You can use a combination of `reshape`, `permute` to map the 3D array to a $n \times 3$ array of pixel values (where n is the number of pixels in the image), and back again**

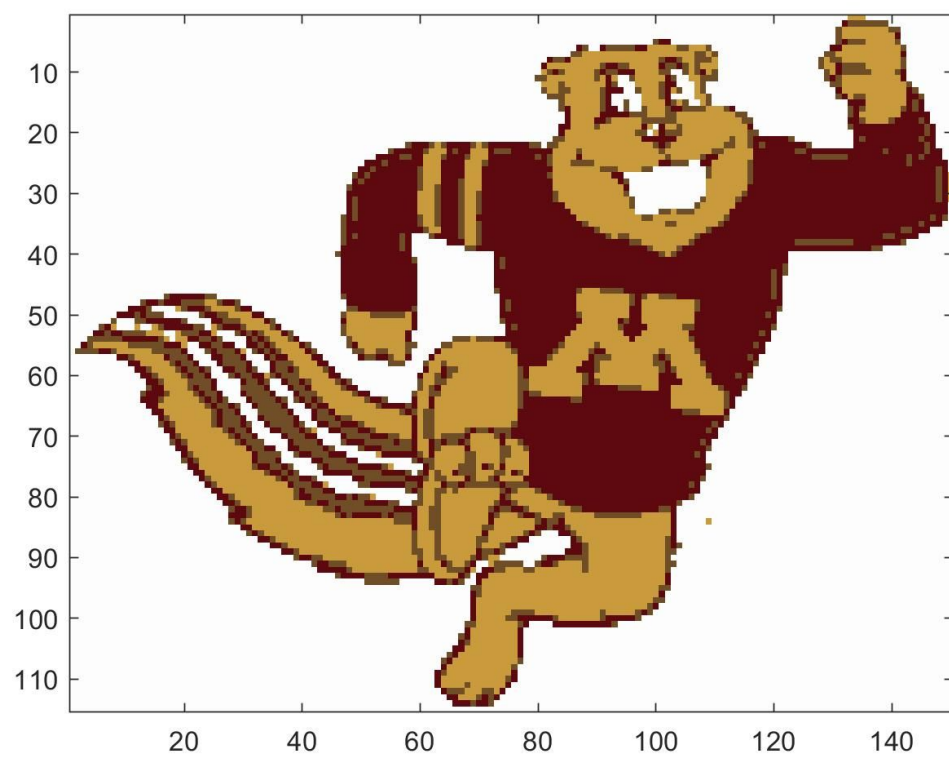
For this part of the Homework I first converted the image matrix (3 Dimensional) Data into a 2Dimensional Feature Matrix where each row represents the RGB Channel values for each pixel. I used the `reshape` command to reshape the $M \times N \times 3$ matrix into a $(M*N) \times 3$ feature matrix. After that I applied the K-Means algorithm on the feature matrix with randomly initialized clusters. After that I take each and every pixel which has been assigned to a particular cluster and change its RGB values based on the RGB values denoted by the cluster center. After this step I again convert the modified feature matrix into a 3 Dimensional matrix which can represent the image with RGB channels. The results is as follows:

Goldy Images:

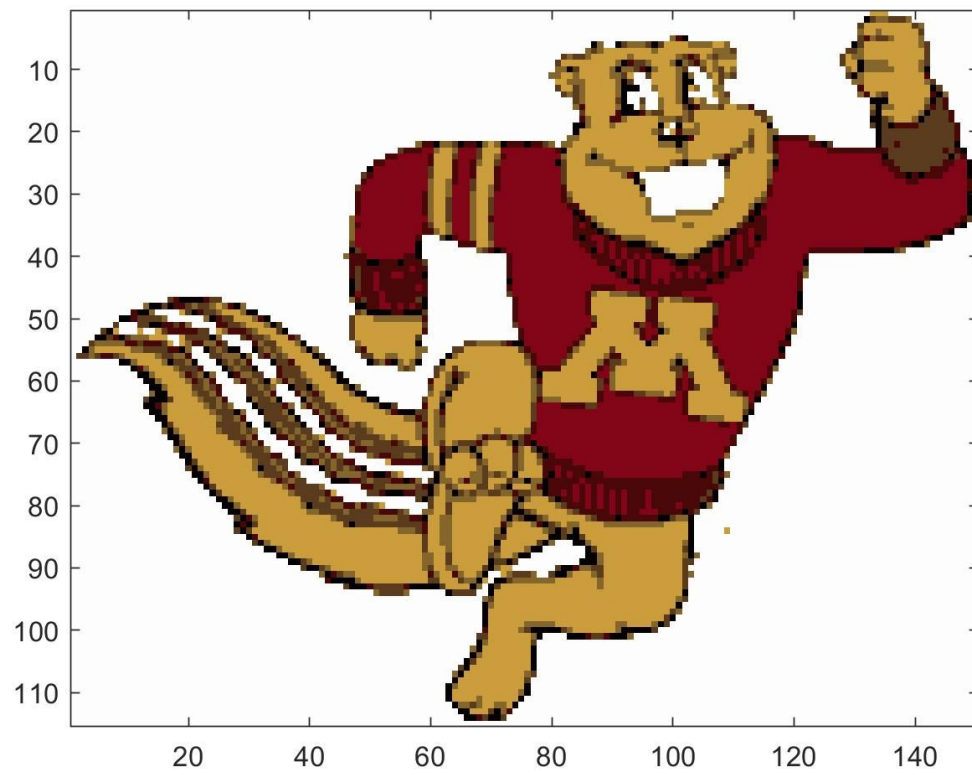
- a. $K = 3$



b. $K = 4$

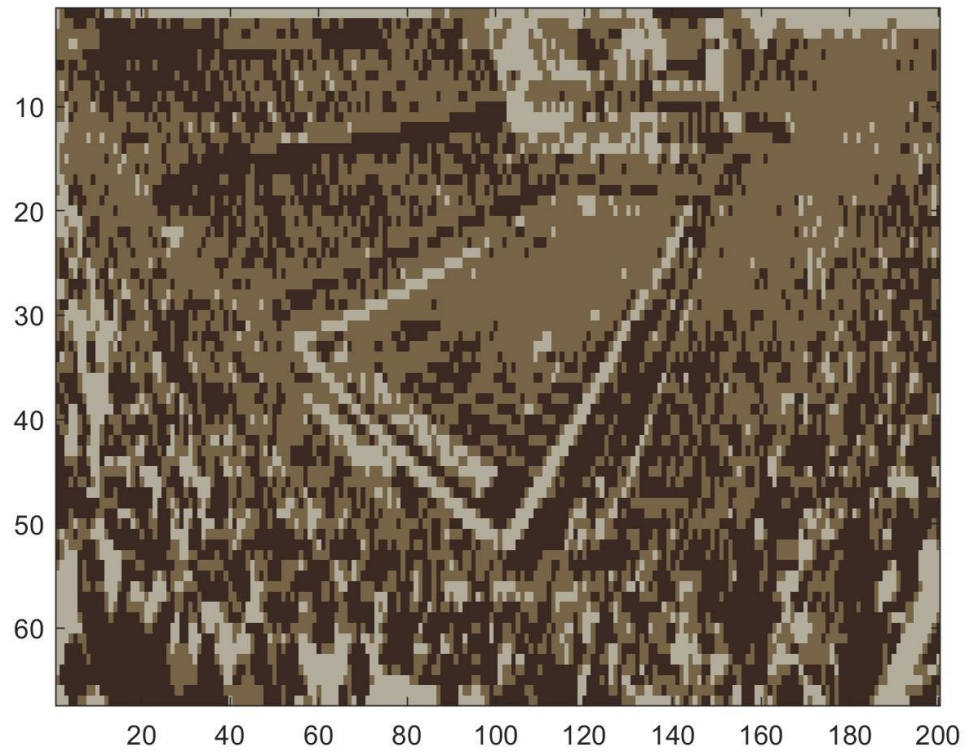


c. $K = 7$

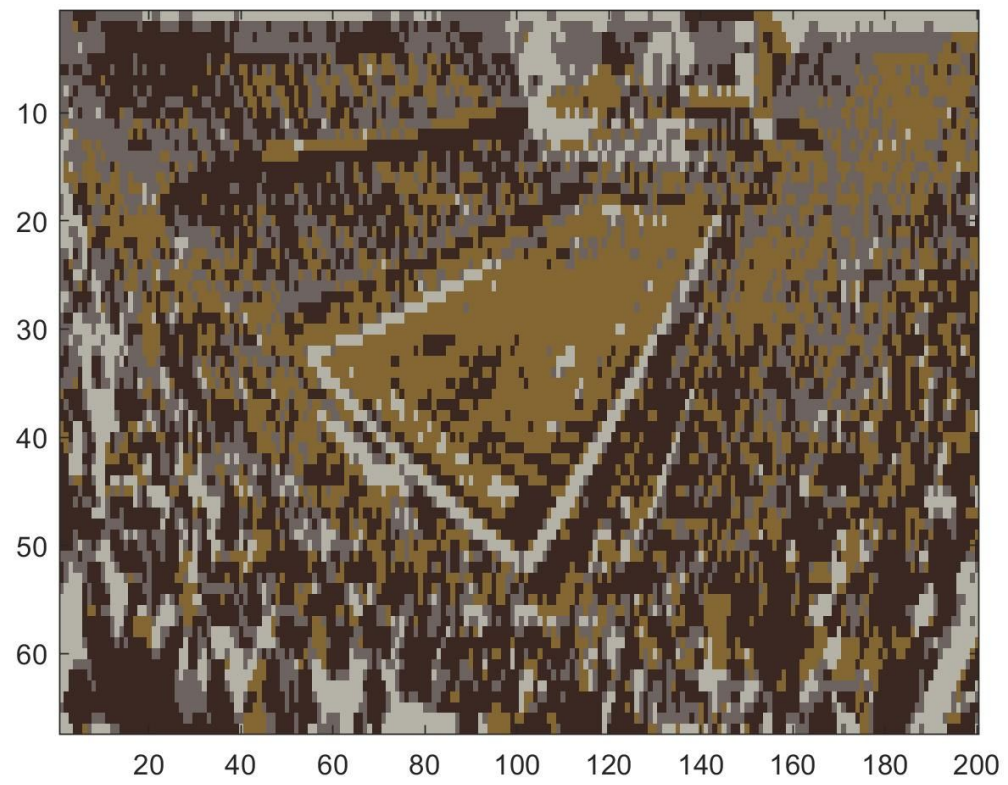


For Stadium

a. $K = 3$



b. $K = 4$



c. $K = 7$

