

Machine Learning Homework 6

Kernel K-means and Spectral Clustering

Due Date 23:55 29th Dec.

- I. **Homework Objective:** Use whatever your favorite language to code out kernel k-means, spectral clustering (both normalize cut and ratio cut). You should consider spatial similarity and color similarity upon the clustering.
- II. **Data:** Two 100*100 images are provided, and each pixel in the image should be treated as a data point, which means there are 10000 data points in each image.
- III. **Kernel:** For both kernel k-means and spectral clustering, please use the new kernel defined below to compute the Gram matrix.

$$k(x, x') = e^{-\gamma_s \|S(x) - S(x')\|^2} \times e^{-\gamma_c \|C(x) - C(x')\|^2}$$

This new defined kernel is basically multiplying two RBF kernels in order to consider spatial similarity and color similarity at the same time. $S(x)$ is the spatial information (i.e. the coordinate of the pixel) of data x , and $C(x)$ is the color information (i.e. the RGB values) of data x . Both γ_s and γ_c are hyper-parameters.

IV. Requirements:

- (20%) You need to make videos or GIF images to show the **clustering procedure** (visualize the cluster assignments of data points in each iteration, colorize each cluster with different colors) of your **kernel k-means**, **spectral clustering** (both normalize cut and ratio cut) program. (Hint : Numpy can help you to solve the eigenvalue problem.)
- (15%) In addition to cluster data into 2 clusters, try more clusters (e.g. 3 or 4) and show your results.
- (15%) For the **initialization of k-means clustering** used in **kernel k-means** and **spectral clustering** (both normalize cut and ratio cut), try different ways and show corresponding results, e.g. k-means++.
- (15%) For **spectral clustering** (both normalize cut and ratio cut), you can see if data points within the same cluster do have the **same** coordinates in the eigenspace of graph **Laplacian**. You should **plot the result and discuss** it in the report.
- (35%) Report: Submit a report in pdf format for showing your code with detailed explanations, giving detailed discussion on experiments as well as your observations. You should explain everything you have done in this homework and show all your results in the report. The report should be written in English.

V. **Turn in:**

1. Report (.pdf)
2. Source code
3. Videos or GIF images of clustering procedure

You should zip all above in one file and name it like ML_HW6_yourstudentID_name.zip, e.g. ML_HW6_0856XXX_王小明.zip.

P.S. If the zip file name has format error or the report is not in pdf format, penalty will be imposed (-10). Please submit your homework before deadline, late submission is not allowed.

◆ Packages allowed in this assignment:

You are only allowed to use numpy, scipy.spatial.distance, package for reading image and visualizing results. Official introductions can be found online.

Important: scikit-learn and SciPy is not allowed.