## Machine Learning: Assignment #5 Fall 2020

**Due:** December 30th, 23:59:59 CST (UTC +8).

## 1. K-Means Clustering

The codes of this section are in the *kmeans* folder.

We will run our first unsupervised algorithm – k-means clustering. Implement k-means algorithm (in kmeans.py), then answer the following questions.

Note that there are different kind of methods to setup initial cluster centers for k-means algorithm, we will use a simple one - randomly choose K samples from dataset as initial cluster centers.

- (a) Run your k-means algorithm on  $kmeans\_data.mat$  with the number of clusters K set to 2. Repeat the experiment 1000 times. Use  $kmeans\_plot.py$  to visualize the process of k-means algorithm for the two trials with largest and smallest SD (sum of distances from each point to its respective centroid).
- (b) You should observe the issue that the outcome of k-means algorithm is very sensitive to cluster centroids initialization from the above experiment. How can we get a stable result using k-means?
- (c) Run your k-means algorithm on the digit dataset  $digit\_data.mat$  with the number of clusters K set to 10, 20 and 50. Visualize the centroids using  $show\_digit.py$ . You should be able to observe that k-means algorithm can discover the patterns in dataset without any label information.
- (d) Another important application of k-means is Vector quantization<sup>1</sup>. Vector quantization is a classical quantization technique from signal processing. It works by dividing a large set of points (vectors) into groups, then representing the data points by their group centroid points, as in k-means and some other clustering algorithms.

Here we will use vector quantization to do image compression. By clustering image pixel value into K groups, we can represent each pixel with log(K) bits, instead of 24 bits (RGB, each channel has 8bit depth).

Finish *vq.ipynb*. Compress images with K set to 8, 16, 32 and 64. I have provided you some sample images, however use your own photos is encouraged.

What is the compress ratio if we set K to 64 (Optionally, you can compute the compress ratio using Huffman encoding)?

https://en.wikipedia.org/wiki/Vector\_quantization

## 2. Spectral Clustering

In this problem, we will try a dimensionality reduction based clustering algorithm – Spectral Clustering. Implement Spectral Clustering algorithm in *spectral.py*, then answer the following questions.

- (a) We will first experiment Spectral Clustering on synthesis data (in *spectral\_exp.ipynb*). Finish *knn\_graph.py* to construct the KNN graph<sup>2</sup> W, then test your algorithm on data *cluster\_data.mat*. Visualize the clustering result and compare it with k-means.
  - Note that we only care about the local connectivity of data points, therefore edges between far away points should be discarded. You should take care with this issue when implementing  $knn\_graph.py$  and choose a proper threshold.
- (b) Now let us try Spectral Clustering on real-world data (in *spectral\_exp.ipynb*). The TDT2 corpus consists of data collected during the first half of 1998 and taken from 6 sources, including 2 newspapers (APW, NYT), 2 radio programs (VOA, PRI) and 2 television programs (CNN, ABC).

Try Spectral Clustering on a subset of TDT2 corpus (in *TDT2\_data.mat*) and compare the result with k-means. You should evaluate the quality of clustering using accuracy and normalized mutual information<sup>3</sup>.

For this part of problem, use *constructW.py* to construct the graph W, and choose proper parameters. As usual, you should preprocessing the data and repeat the experiments multiple times then take the average.

## 3. Principal Component Analysis

Principal component analysis (PCA) is a mathematical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components. Let us deepen our understanding of PCA (you should implement it in pca.py) by the following problems.

- (a) Our actions to hack the CAPTCHA system last time have been detected, so they updated the system to add rotation in the CAPTCHA images:
  - Luckily, we have just learnt PCA. Your task is to implement *hack\_pca.py* to recover the rotated CAPTCHA image using PCA.
  - Some rotated CAPTCHA samples are provided for you. It's OK for your algorithm to return grayscale image or image with size different form input image.
- (b) Now let us apply PCA to a face image dataset in ORL\_data.mat.

<sup>&</sup>lt;sup>2</sup>Connect each data point with its K nearest neighbors, you can use binary edge weights or compute edge weights using Heat kernel.

<sup>&</sup>lt;sup>3</sup>See http://www.cad.zju.edu.cn/home/dengcai/Data/Clustering.html and Deng Cai, Xiaofei He, and Jiawei Han, "Document Clustering Using Locality Preserving Indexing", in IEEE TKDE, 2005, for more details.



- (i) Run PCA on the dataset, visualize the learnt eigenvectors using *show\_face.py*. You should see faces in the image, these faces are called Eigenface.
- (ii) Use PCA to do dimensionality reduction<sup>4</sup>. You should experiment with number of reduced dimensionality being 8, 16, 32, 64 and 128 respectively. For each test, run KNN on the dimensionality reduced dataset and report the testing error rate.
- (iii) Will dimensionality reduction cause loss of information? Let us see it visually. Do dimensionality reduction using PCA then recover the original image form the low dimensional feature vectors. Visualize the original and recovered image using show\_face.py. Again, you should experiment with number of reduced dimensionality being 8, 16, 32, 64 and 128 respectively.

Please submit your homework report to at http://courses.zju.edu.cn:8060/course/18843 in pdf format, with all your code in a zip archive.

<sup>&</sup>lt;sup>4</sup>You are strongly encouraged to try LDA on this part of problem. You should be impressed by its outstanding performance.