

CS5487 Machine Learning - Programming Assignment 2

Zhanghan Ke (55460880)

zhanghake2-c@my.cityu.edu.hk

Abstract. This is the report for my Programming Assignment 2 of CS5487 Machine Learning. The code also available in my github. Link: <https://github.com/ZHKKe/PA2.git>.

1 Clustering Synthetic Data

1.1 Problem (a)

I implemented K-Mean and Mean-Shift algorithms, and the EM-GMM algorithm of my assignment from the *scikit learn* package.

1.2 Problem (b)

The Fig. 1 shows my results of this problem. The results shows the Mean-Shift with careful selection of bandwidth h could clustering data best. And the EM-GMM algorithm also work pretty good. The K-Mean gains a bad result in the *dataB_X*.

K-Mean assigns each sample to a particular cluster on convergence. It works better in the high dimensional data nad it is easy to interpretation and implementation. However, K-Mean might cause some samples mis-matching to any cluster. EM-GMM assigns each sample to clusters softly by a probability, it works well with non-linear geometric distribution as well. But since the all components must be accessed, it is difficult in the high dimensions. Mean-Shift has a hyperparameter bandwidth h , which could control the clustering results, it could be set case by case to gain great results. But this algorithm is very sensitive to h , so we need to choose it carefully.

1.3 Problem (c)

Mean-Shift is sensitive for the bandwidth, changing it will cause bad results. The Fig. 2 shows the result of *dataC_X* under different h .

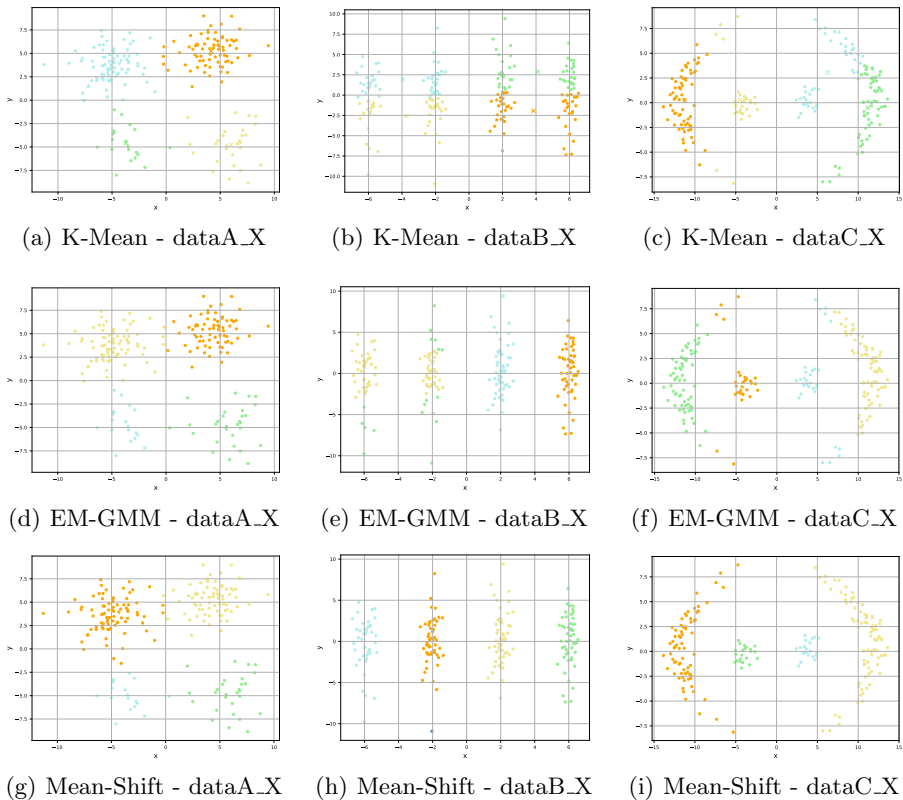


Fig. 1. Result of clustering synthetic data by three algorithm.

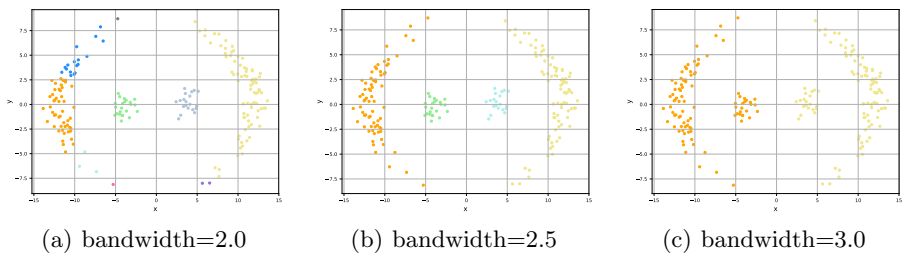


Fig. 2. Result of clustering synthetic data of Mean-Shift under different bandwidth. (a) Too small bandwidth will generate too many clusters. (b) Suitable bandwidth will generate matching clusters. (c) Too large bandwidth will generate too less clusters.

2 A Real World Clustering Problem Image Segmentation

2.1 Problem (a)

The Fig. 3 shows the results in 2 different images by three algorithm.

2.2 Problem (b)

I try to use different setting of the distance and kernel and find these parts need to be setting carefully, otherwise the results will be bad.

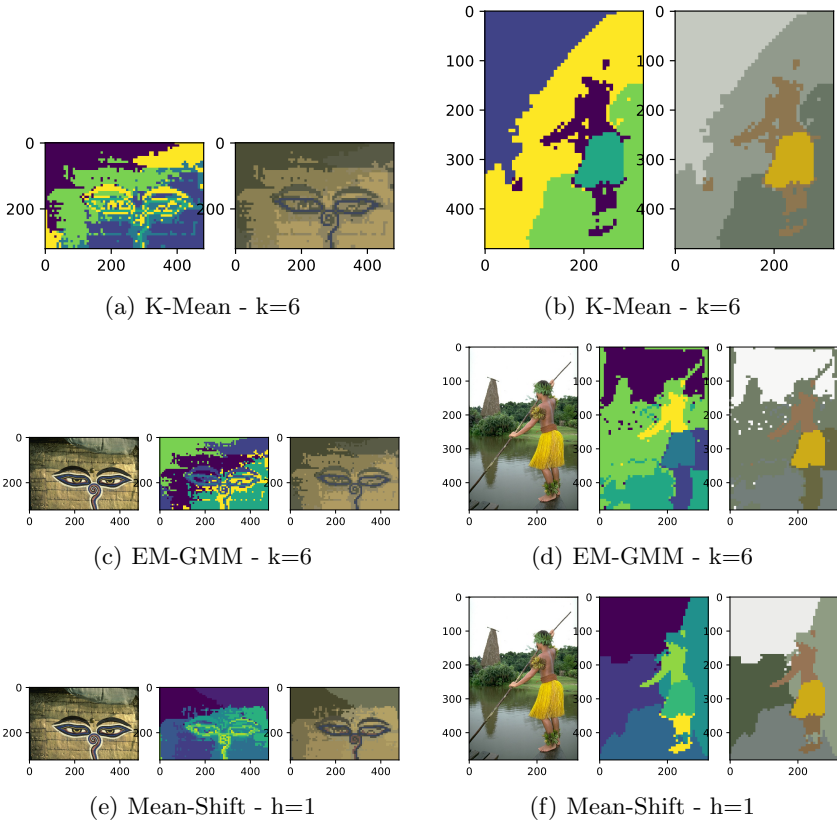


Fig. 3. Result of image segmentation by three algorithm.