

EE 219 Project 2

Clustering

Winter 2018

Xiaohan Wang

Yang Tang

Introduction

In this project, we did clustering algorithms on textual dataset “20 Newsgroup” to capture the dependencies among the documents and find proper representations of the data pretending not knowing the class labels of the documents. Then we use the class labels as ground truth to evaluate the performance of the clustering task. Besides, we tried different preprocess methods to improve the performance of the clustering algorithms.

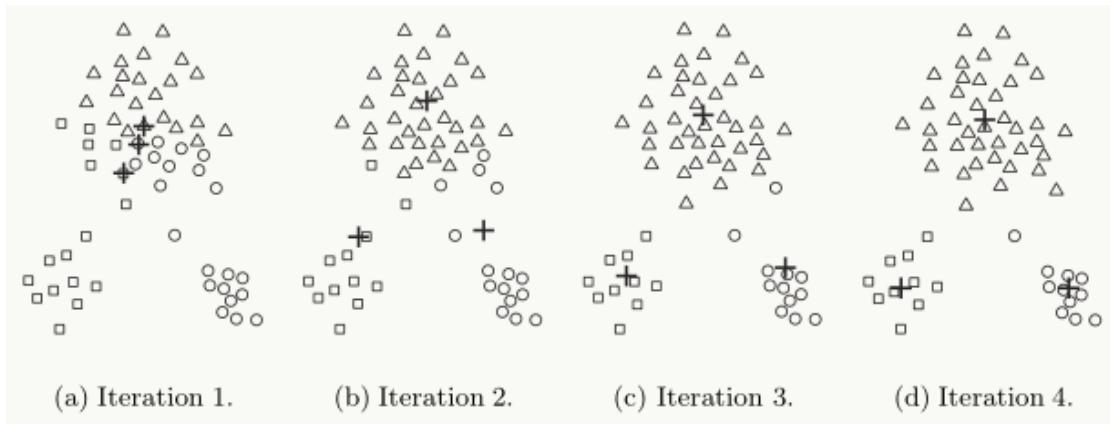
In project 1, we performed classification on the same dataset. These two methods, classification and clustering, though, are not the same cases. Their goals are similar in that they both need to categorize the data. However, classification is a priori algorithm. We have some given categories first, then we need to train and extract the feature of those categories, and finally utilize those features to classify the testing data into proper given categories. Clustering makes a different story. Clustering algorithms are unsupervised methods for finding groups of data points that have similar representations in a proper space. It has no priori labeling(grouping) and we need to predict, test and correct the categories iteratively until all testing data are clustered into proper clusters in convergence.

In this project, we use K-means clustering algorithm. Its idea is similar to least squares method. We performed the following steps to minimize the least squares of distances between each data point and the center of the cluster

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

Where μ_k is the kth cluster center, and r_{nk} denotes the nth data point is belonging to the kth cluster or not.

- 1) Initialize randomly k cluster center positions.
- 2) Assign each data point to the nearest cluster center.
- 3) Calculate the position of each center of cluster.
- 4) Set the center of the cluster to the mean of the data points that are currently within the cluster.
- 5) Iteratively repeat step 2) to 4) until all data points converge to kth cluster centers.



Question 1

Similar to that of project 1, using the same 8 subclasses datasets as table 1, first we extract key terms of documents and quantify how significant a word is to a document by using TFxIDF.

Table 1. Data set classes

Class 1	Class 2
comp.graphics	rec.autos
comp.os.ms-windows.misc	rec.motorcycles
comp.sys.ibm.pc.hardware	rec.sport.baseball
comp.sys.mac.hardware	rec.sport.hockey

The data process is as follow:

1. Remove the punctuations in original data. Punctuations are meaningless to the content of documents, while they can add much noise to feature extraction, so we need to remove those punctuations. We listed all possible punctuations referring to ASCII codes, and then split string using regex tokenizer.
2. Remove newsgroup headers, footer signatures and quotes. Those are irrelevant information to key meaning of the documents.
3. Exclude stop words. Stop words like “I”, “she”, etc. can appear in all texts and cannot contribute the feature to classify the text, so we need to remove them.
4. Create TFxIDF vector to get feature terms. By setting the `min_df` parameter, we can set the baseline of appearance that a candidate word need to have. If a word appears less time than `min_df`, we can ignore it as irrelevant word. The larger the `min_df` is, the less number of featured terms we will get.

Using all sets (train sets and test sets) data of 8 topics, the dimensions of the TFxIDF matrix is 7882×15464 , which means there are 7882 documents with 15464 featured terms extracted.

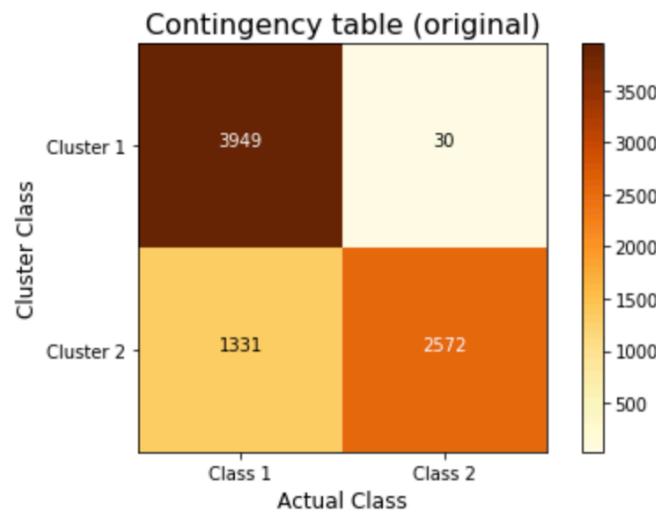
Question 2

In this part, we implemented K-means clustering method setting $k = 2$. Then we visualized our clustering results using contingency table and evaluated the effects with five factors.

Question a)

In project 1, we used confusion matrix to visualize and evaluate our classification results. In this part, to visualize our clustering results, we need to use contingency table which is similar to confusion matrix. Confusion matrix is used to evaluate supervised data with two factors “actual classes” and “predicted classes”. In contingency table, each element A_{ij} denotes the number of data points that are member of class c_i and elements of cluster k_j . Here, c_i can be viewed as “actual classes” and k_j can be viewed as “predicted classes” in confusion matrix. But the difference is that we don’t

have given “actual classes” labels and all we need to do is to categorize the data into several clusters with no care of the clusters’ name. So setting with a 2x2 confusion matrix, we may get a diagonal-matrix or a vice-diagonal-matrix. Then we can adjust the matrix by finding a row permutation to get a diagonal-matrix. Below are our contingency table (diagonal-confusion matrix):



From the table above, we can see that through K-means clustering algorithm, it can hardly get to a diagonal matrix, which indicates that the clustering algorithm is not that ideal to purely cluster the data.

Question b)

In order to make a concrete comparison of different clustering results, we evaluated them through 5 measures: homogeneity score, completeness score, V-measure, adjusted rand score and adjusted mutual info score. Below are our results for K-means clustering results:

Table 2. 5 measures for K-means clustering results

Measures	Values
Homogeneity score	0.4243358018332555
Completeness score	0.4636994594803498
V-measure	0.4431452004668879
Adjusted rand score	0.428504440021173
Adjusted mutual info score	0.4242830978688626

From the table above, we can see that the 5 measures score are distributed in the range between 0.42 and 0.47 with the best value should be 1. The more the contingency table is like a diagonal matrix, the better our clustering algorithm clustered the data points purely and correctly. Clearly, our clustering algorithm now is not as good as possible, and it's not very stable with some fluctuations in each time's results.

Question 3

Question i

From the results in question 2, we can see that though the contingency table seems like a diagonal matrix, its performances is still not ideal. And this may be caused by the high dimension of TFxIDF matrix. Thus, in the next step, we need to reduce its dimension and find the effective dimension of the data through inspection of the top singular values, i.e. choosing the number of principal components.

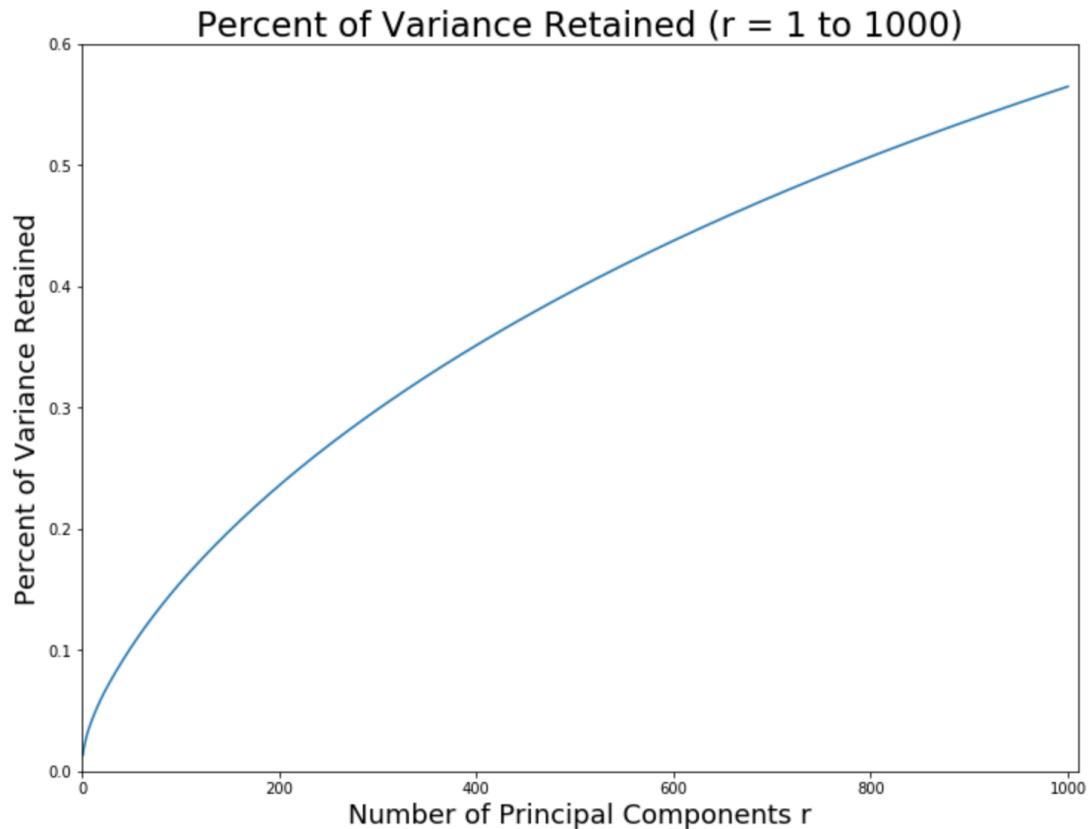
There is an equation in machine learning to evaluate the effects of dimension r

$$err = \frac{\frac{1}{m} \sum_1^m \|x_i - x_{i,approx}\|^2}{\frac{1}{m} \sum_1^m \|x_i\|^2}$$

if err is less than 0.01, than we can say that “99% of variance is retained”. First we set $r = 1$, and calculate its err and check if it can satisfy the requirement. If not, then we increase r by 1 each time until finding the satisfied value. Here we can use truncated SVD to easily calculate err by

$$err = \frac{\sum_1^r S_{ii}}{\sum_1^n S_{ii}}$$

Below are our results with having r from 1 to 1000:



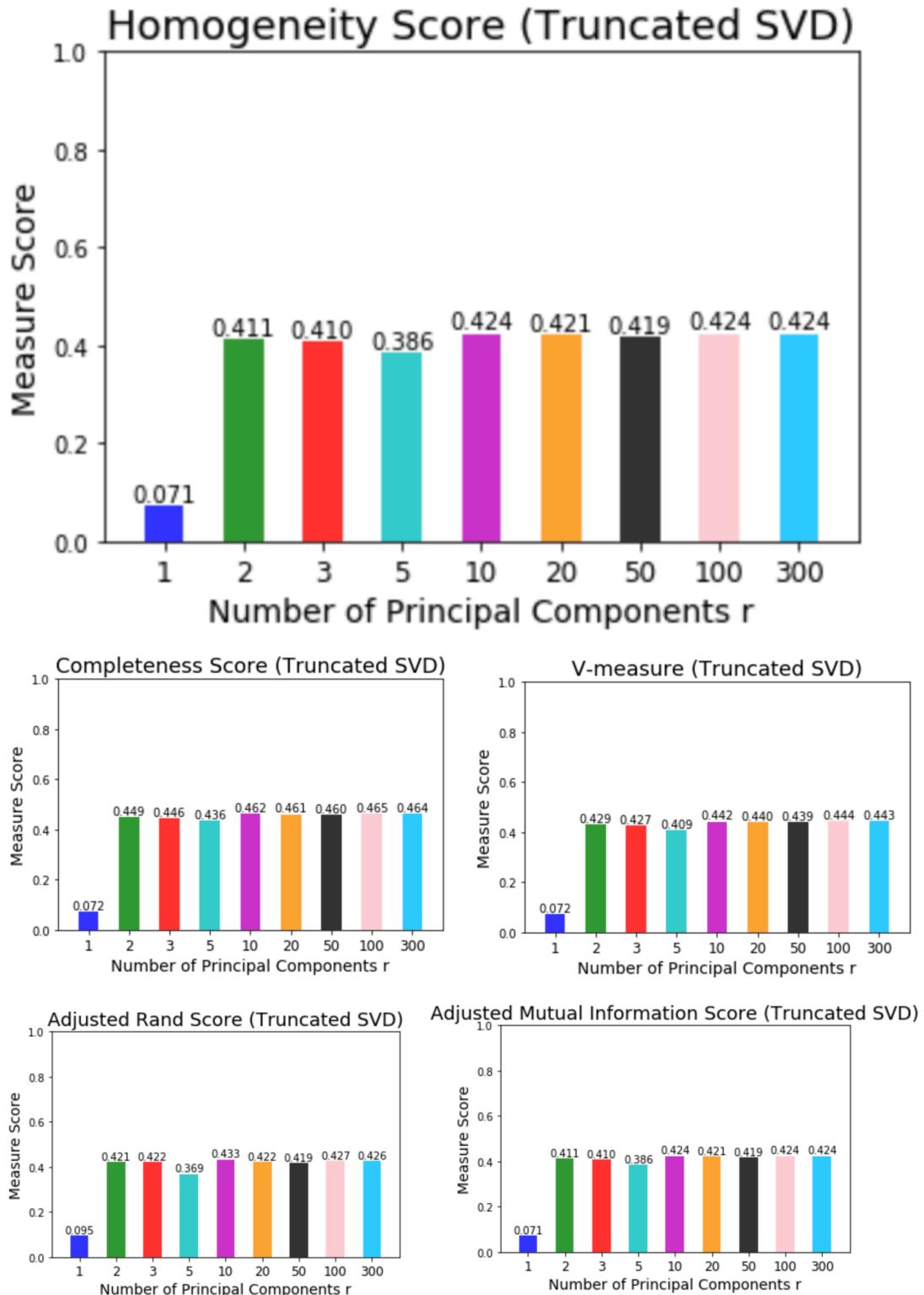
From the figure above, we find that as r increases, the percentage of variance the top r principal components can retain increases monotonically.

Question ii

In this part, we use two methods to reduce dimension and compare their effects. One is Latent Semantic Indexing (LSI), the other is Non-Negative Matrix Factorization (NMF).

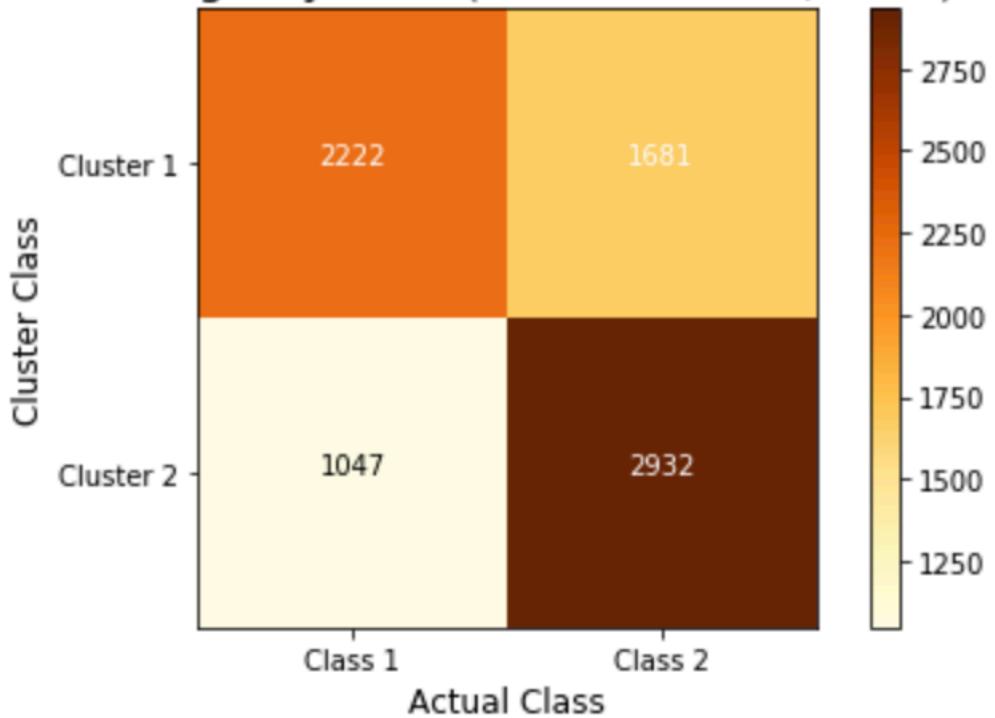
1. Truncated SVD (LSI) / PCA

1) 5 measures scores v.s. r

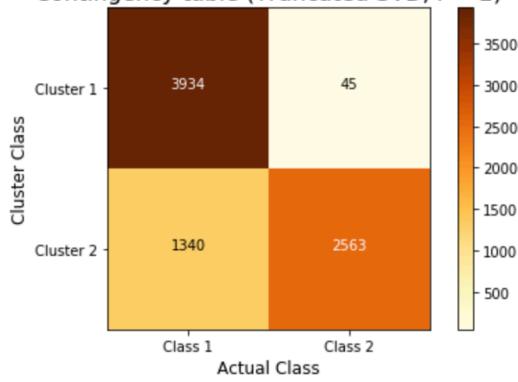


2) Contingency matrices

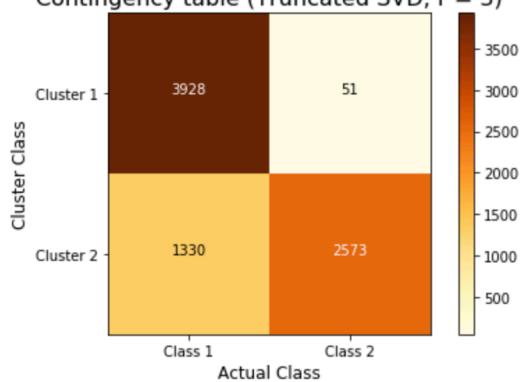
Contingency table (Truncated SVD, $r = 1$)



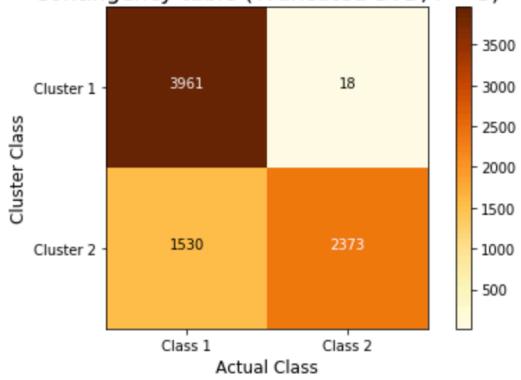
Contingency table (Truncated SVD, $r = 2$)



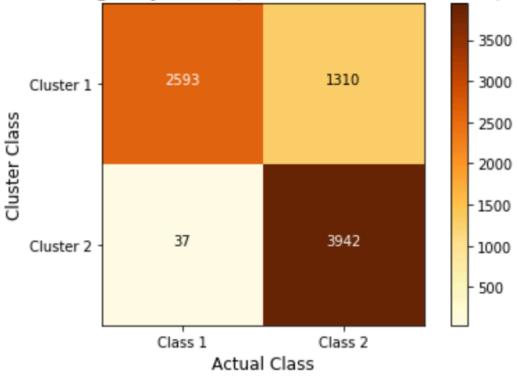
Contingency table (Truncated SVD, $r = 3$)

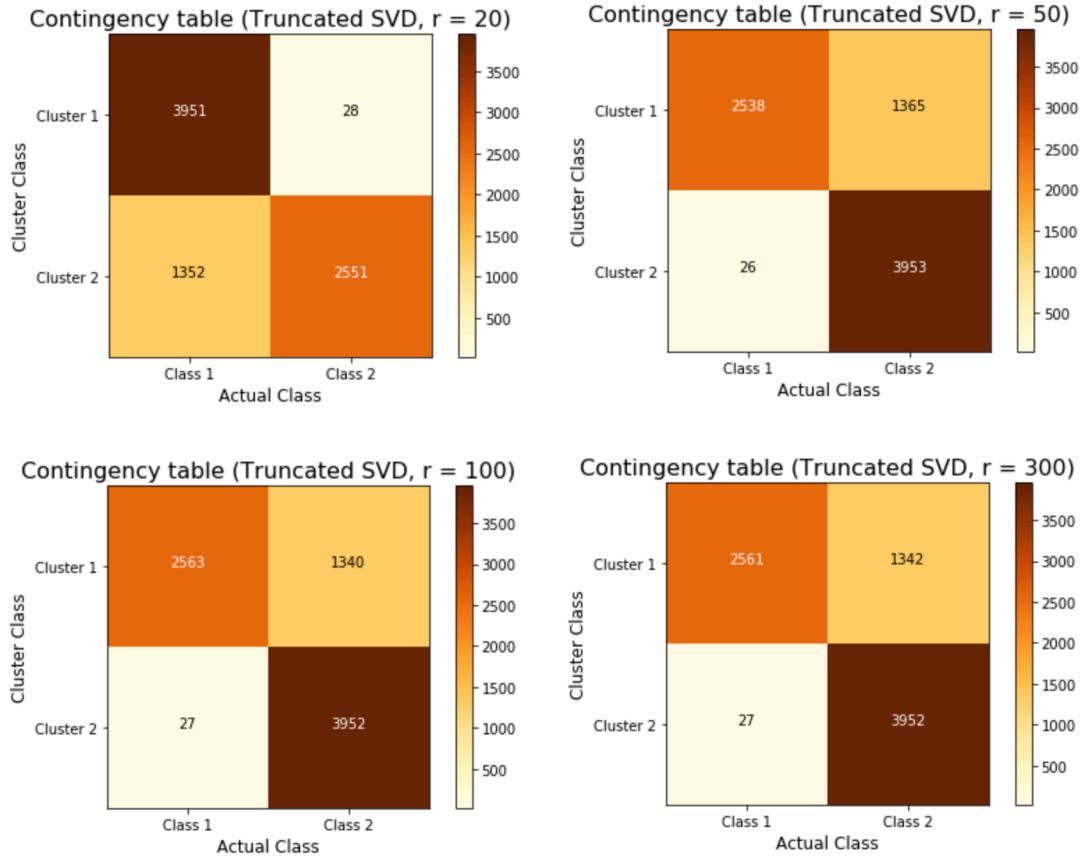


Contingency table (Truncated SVD, $r = 5$)



Contingency table (Truncated SVD, $r = 10$)

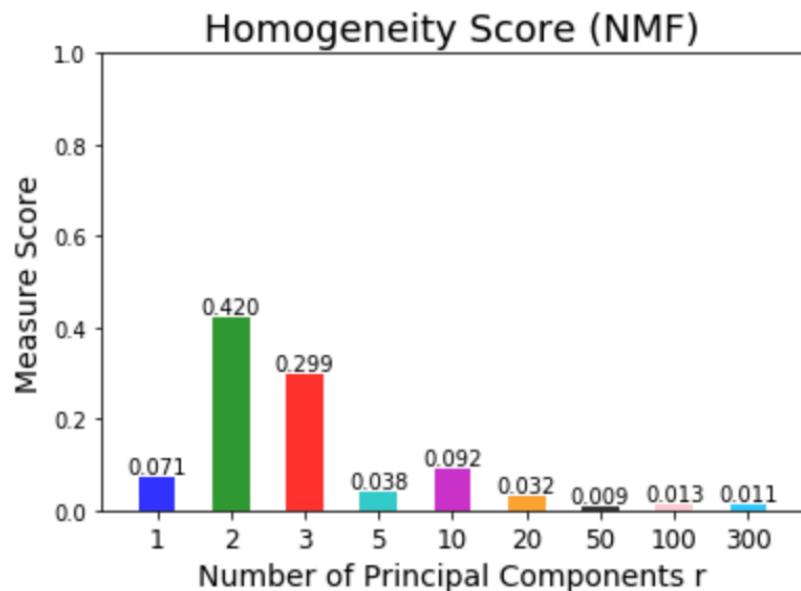


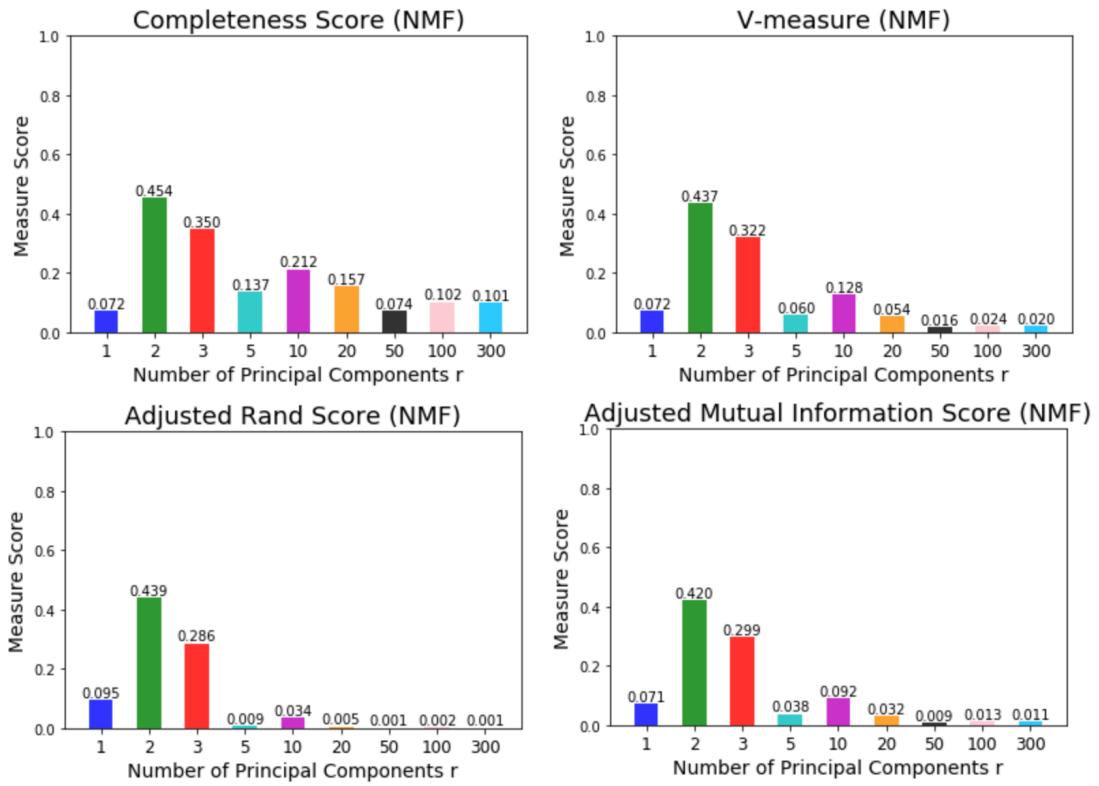


Comparing 5 measures and contingency tables comprehensively, we find that the best r choice for Truncated SVD is 10.

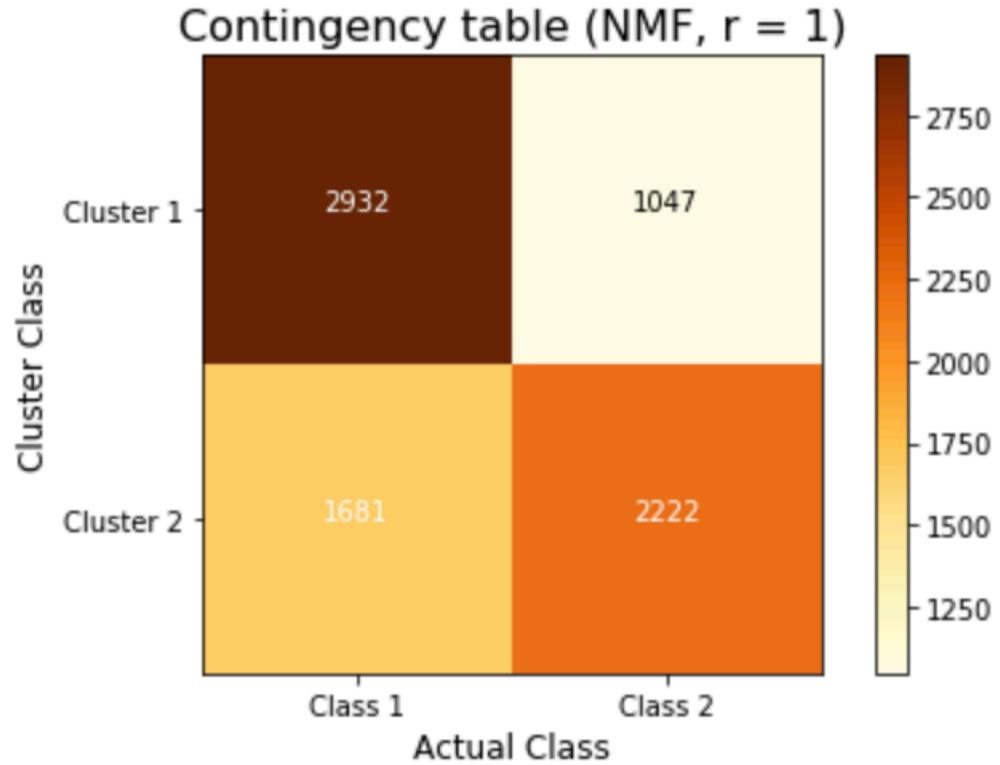
2. Non-Negative Matrix Factorization (NMF)

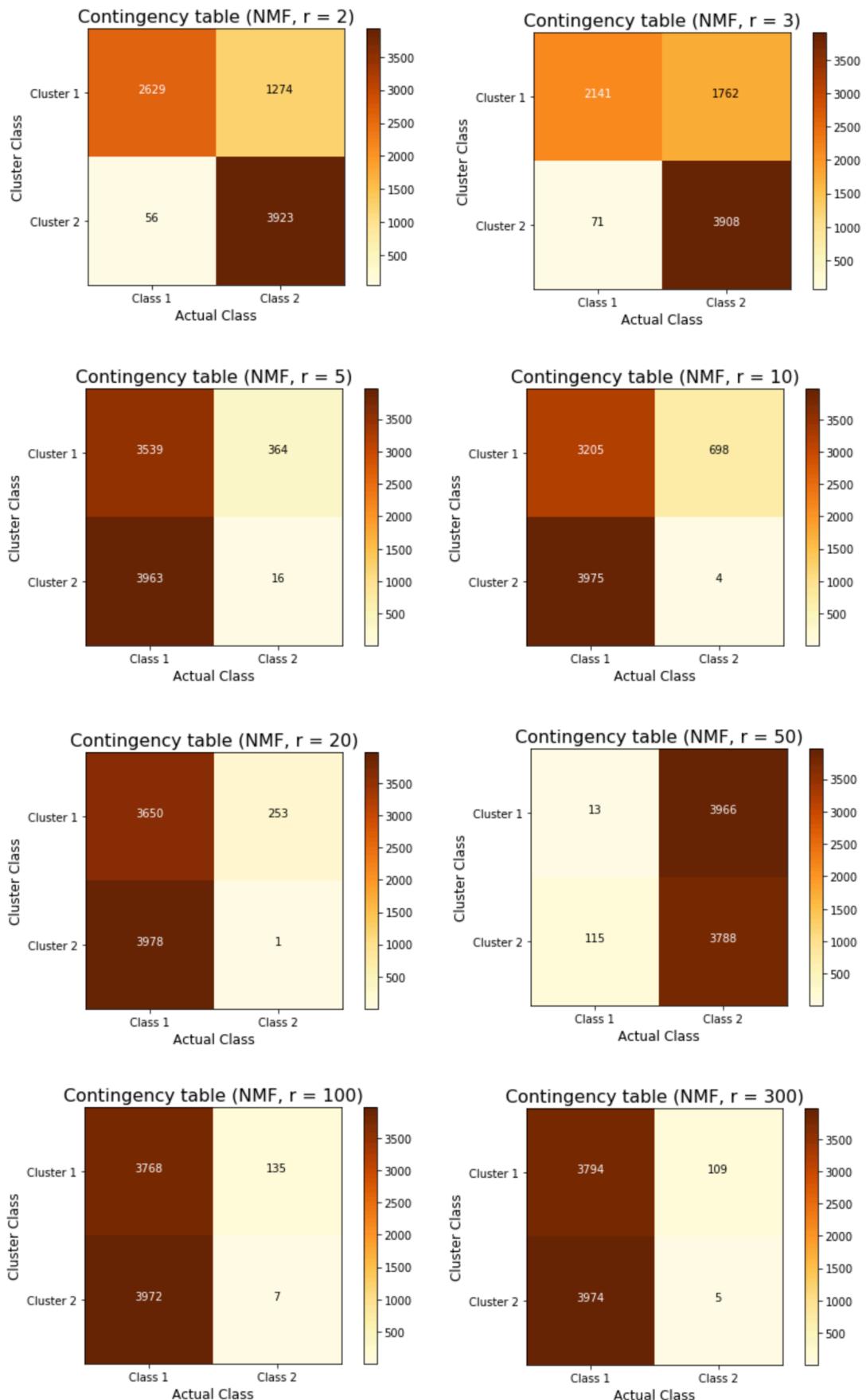
1) 5 measures scores v.s. r





2) Contingency matrices





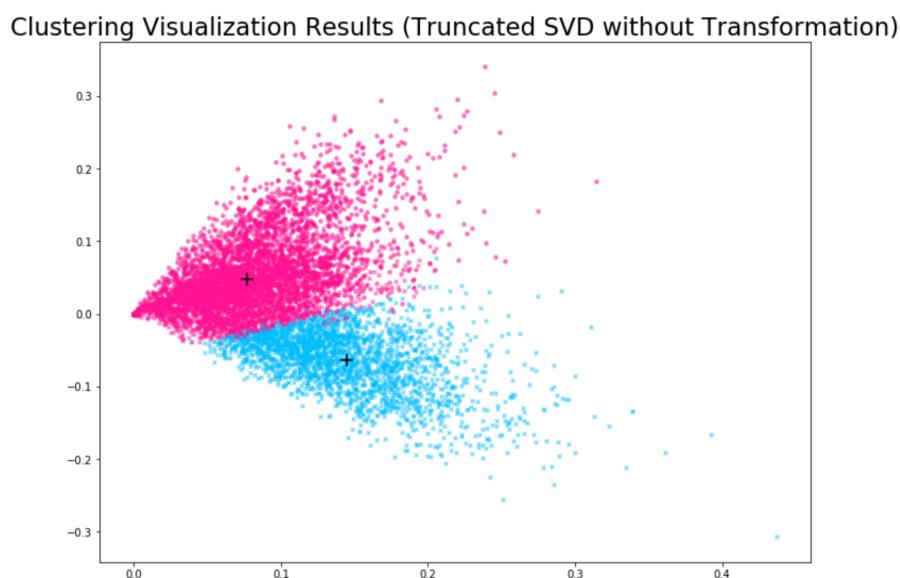
Comparing the 5 measure scores and contingency tables above, we can easily conclude that the best r choice for NMF is 2.

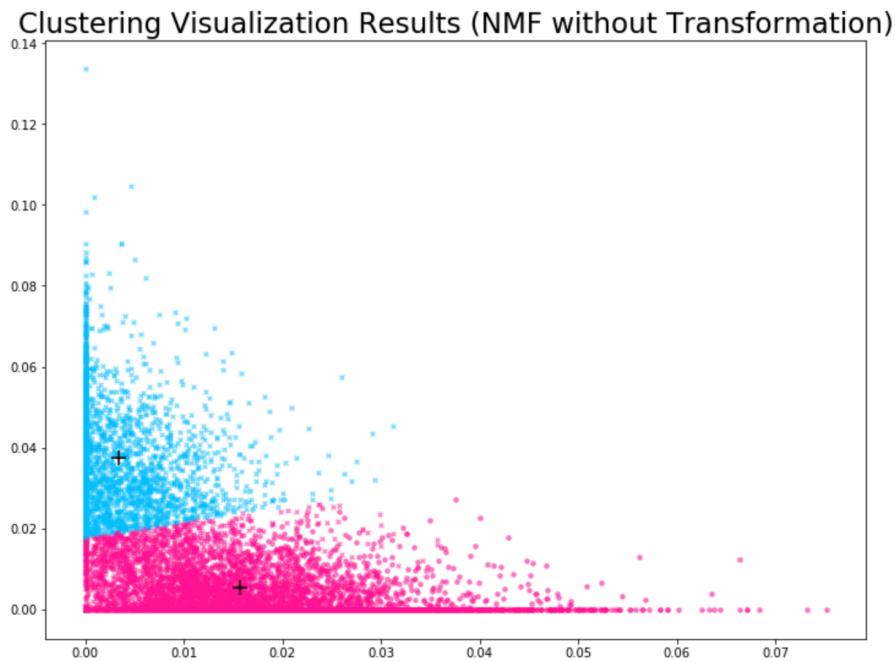
By observing the figures of 5 measures scores, we find that they are non-monotonic as r increases. This is because when r is at too low values, the information that the singular value contains is limited so that we cannot guarantee the correctness of the clustering, which caused the measure scores low at first. As r increases, it contains more information to help clustering accurately, which contributes to the increment of the measure scores. While when r raised up with too big value, it will create a high dimensional matrix. This brings the same problem as in question 2 that the Euclidean distance is not a good metric anymore, in the sense that the distances between data points tends to be almost the same. This leads to the decreasing clustering performance evaluated by 5 measures. Therefore, the 5 measure scores figures are not monotonic, and we need to find the best dimension r in different cases.

Question 4

Question a)

According to the results in question 3, we set the reduction dimension values as: Truncated SVD with $r = 10$, NMF with $r = 2$. Below are our clustering visualization results:

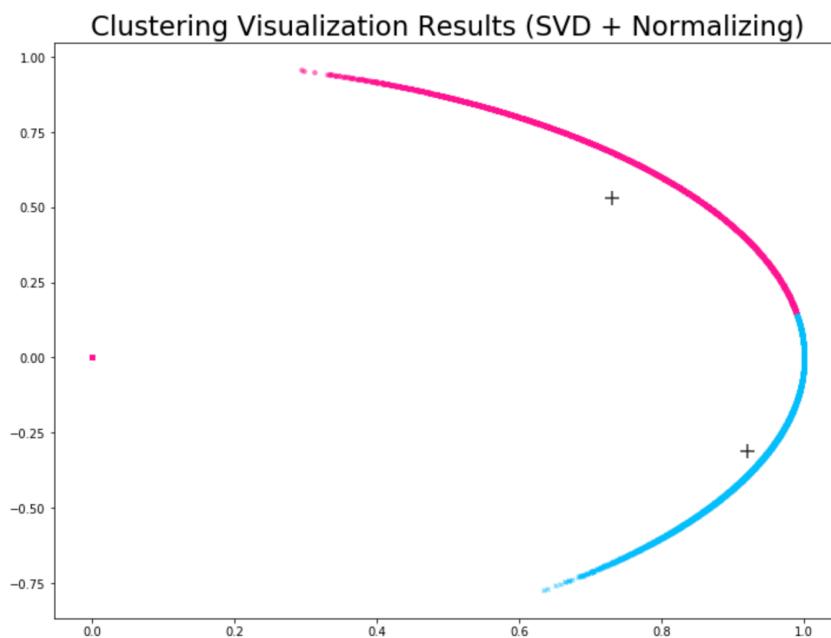


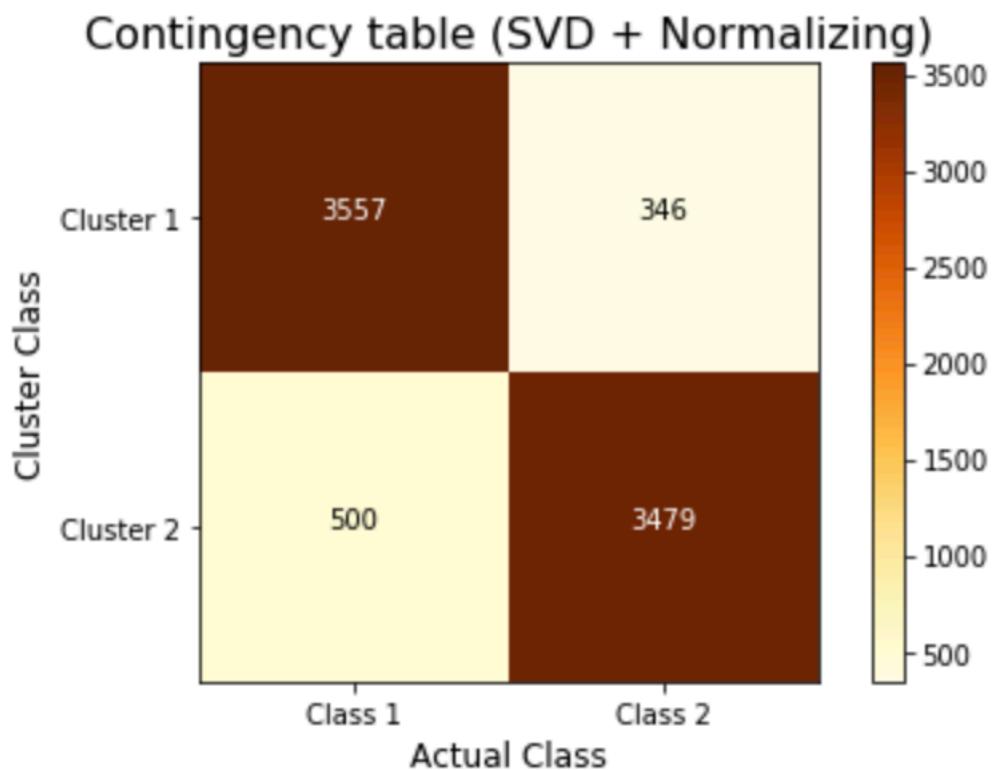


Question b)

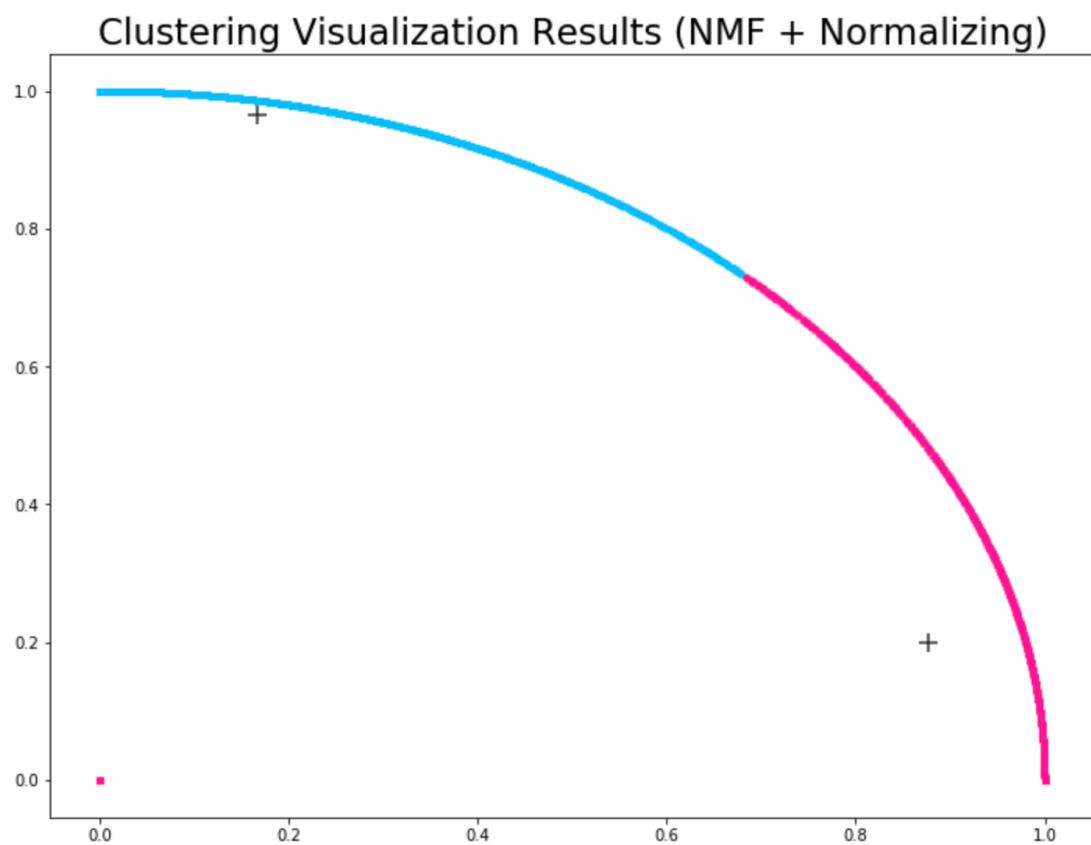
In this part, we tried three methods to improve clustering performance in question a. Same as that in last part, we set the reduction dimension values as Truncated SVD with $r = 10$, NMF with $r = 2$. Below are our clustering visualization results and corresponding contingency tables:

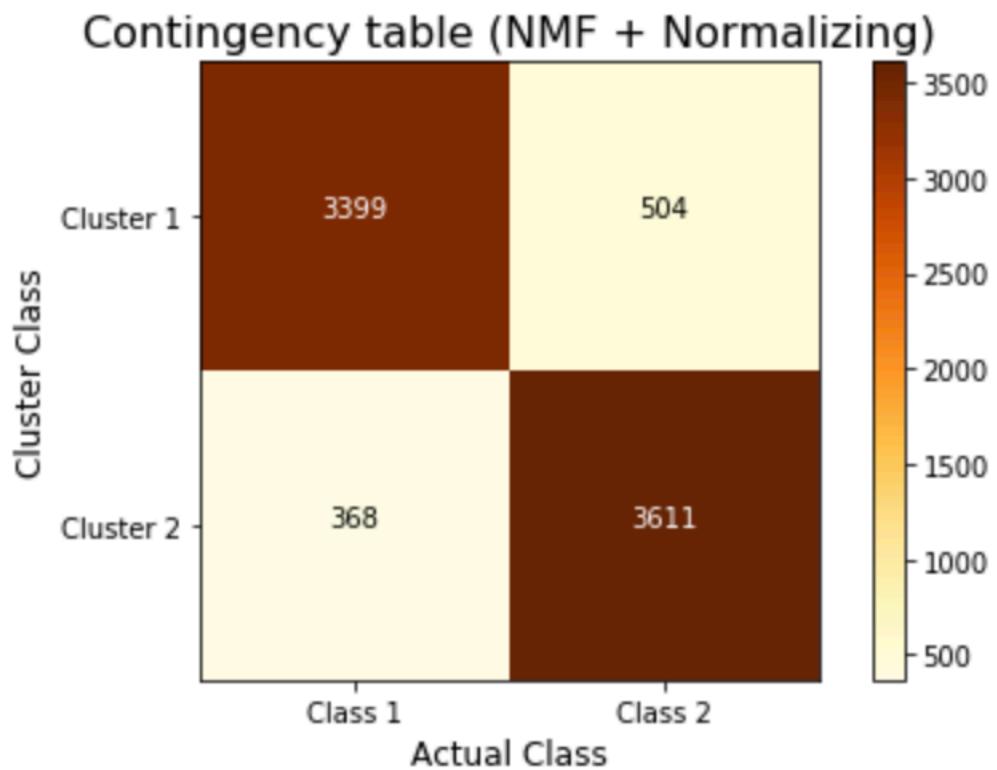
- 1) Truncated SVD + Normalizing



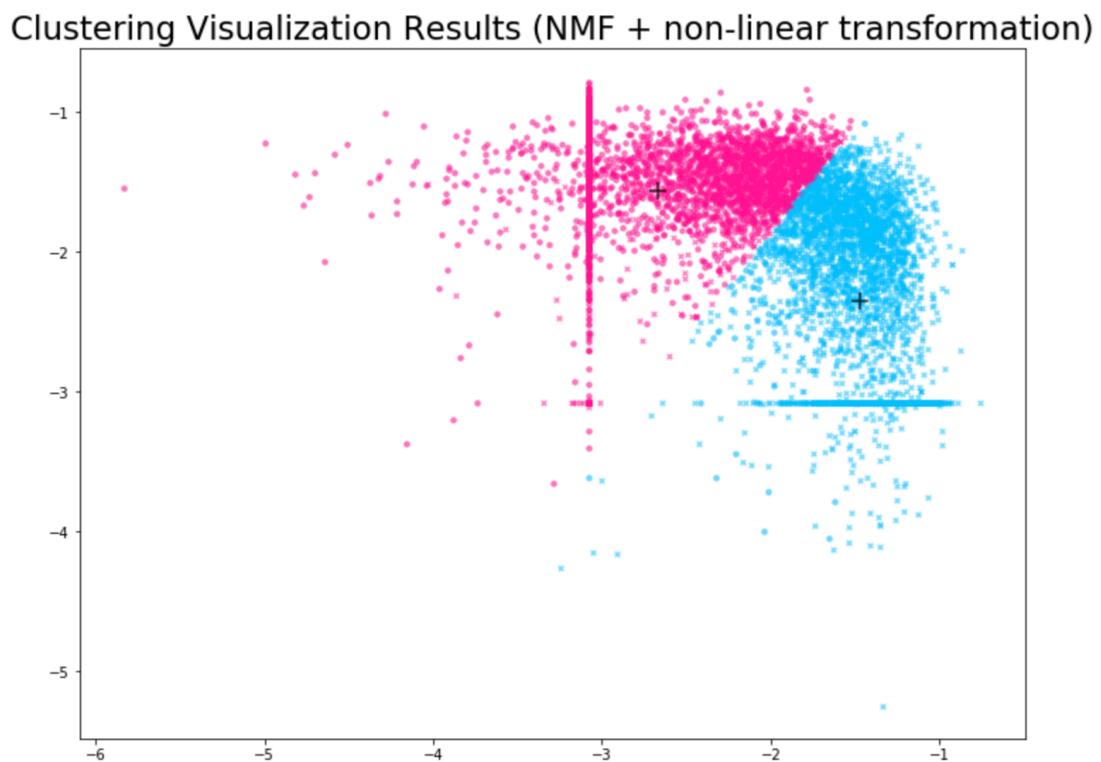


2) NMF + Normalizing

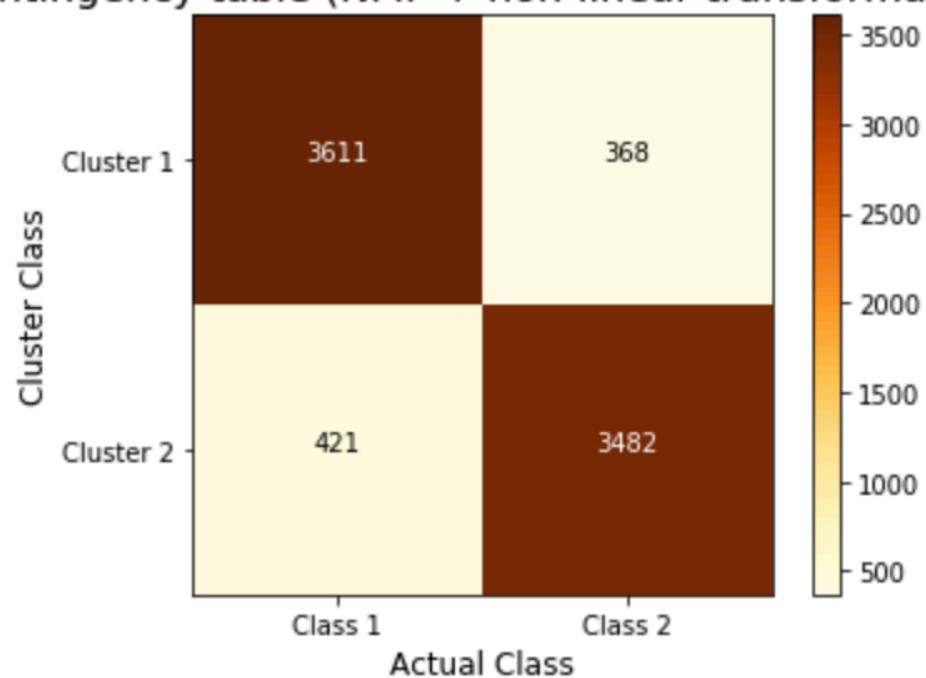




3) NMF + Non-linear Transformation

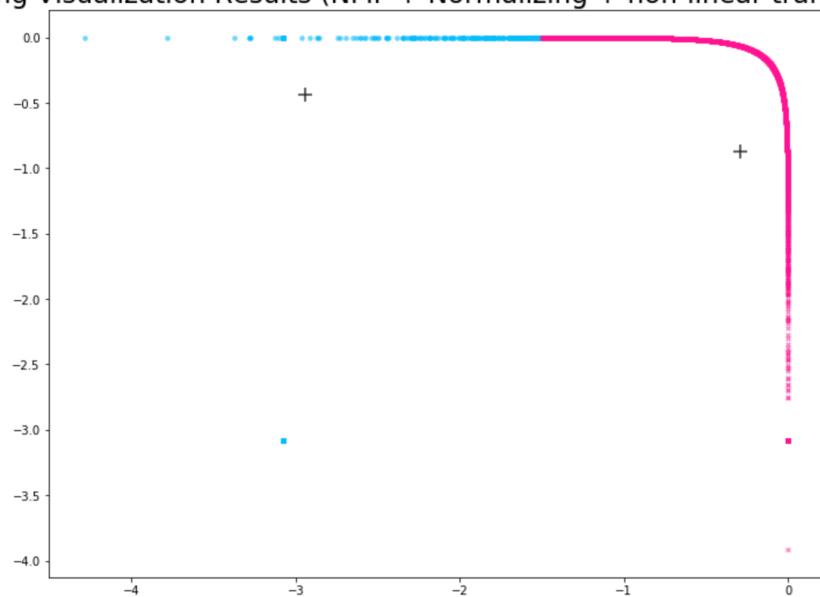


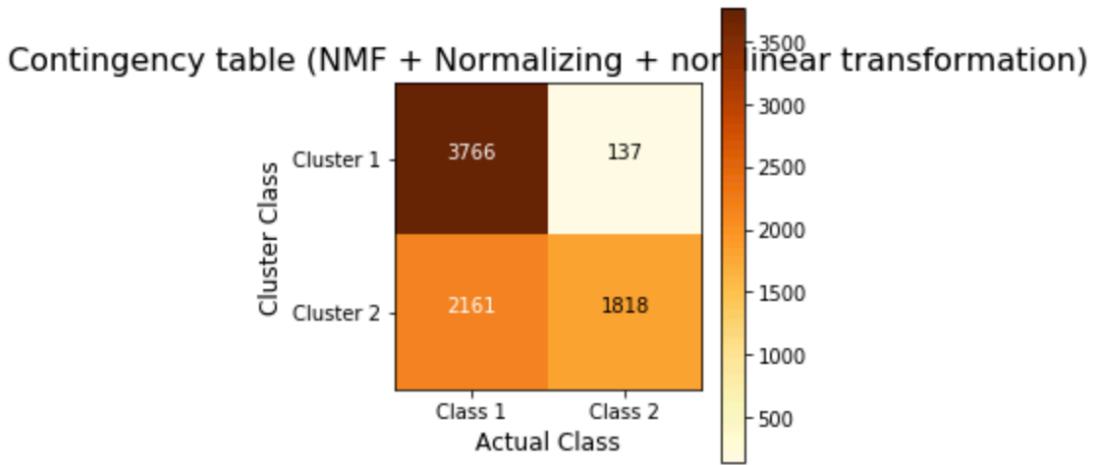
Contingency table (NMF + non-linear transformation)



4) NMF + Normalizing + Non-linear Transformation

Clustering Visualization Results (NMF + Normalizing + non-linear transformation)





From the figures above, we can see that with non-linear (logarithm) transformation, the clustering performance is improved greatly. With only applying normalizing, we can also improve the performance. While if we combine both transformation methods with NMF, the clustering results appeared not ideal as we expect.

The reason why logarithm transformation may increase the clustering results is that the relevance of a term or a document is not proportional to term frequency. With sub-linear logarithm method can help eliminate the effect so that the influence of very large or very small values is amortized, which can contribute to a better performance of clustering.

Question 5

In the final part of project, we did the same dimension reduction then K-means clustering procedure to cluster the documents from 20 topics into 6 clusters (as indicated in table 3) based on the classes they belong to. We assigned label 0 to label 5 to represent actual clusters of original data. And we applied different dimension reduction methods which had been improved in the last part of the project on the dataset before K-means clustering. Also, we chose various numbers of terms to represent each document.

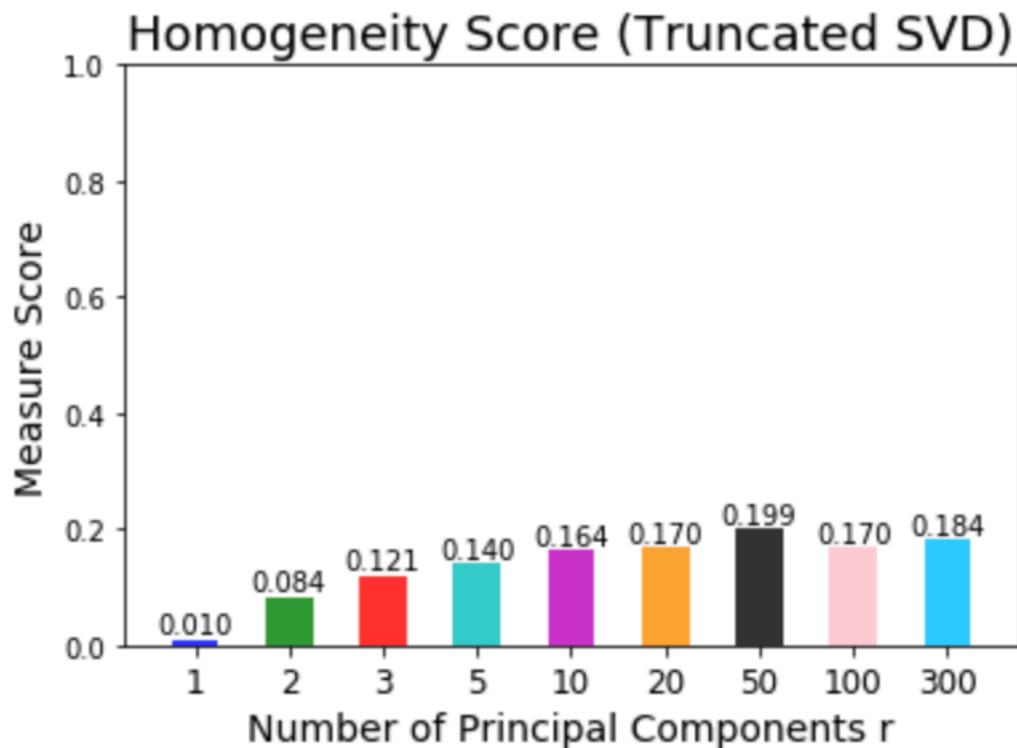
Table 3. All data sets classes

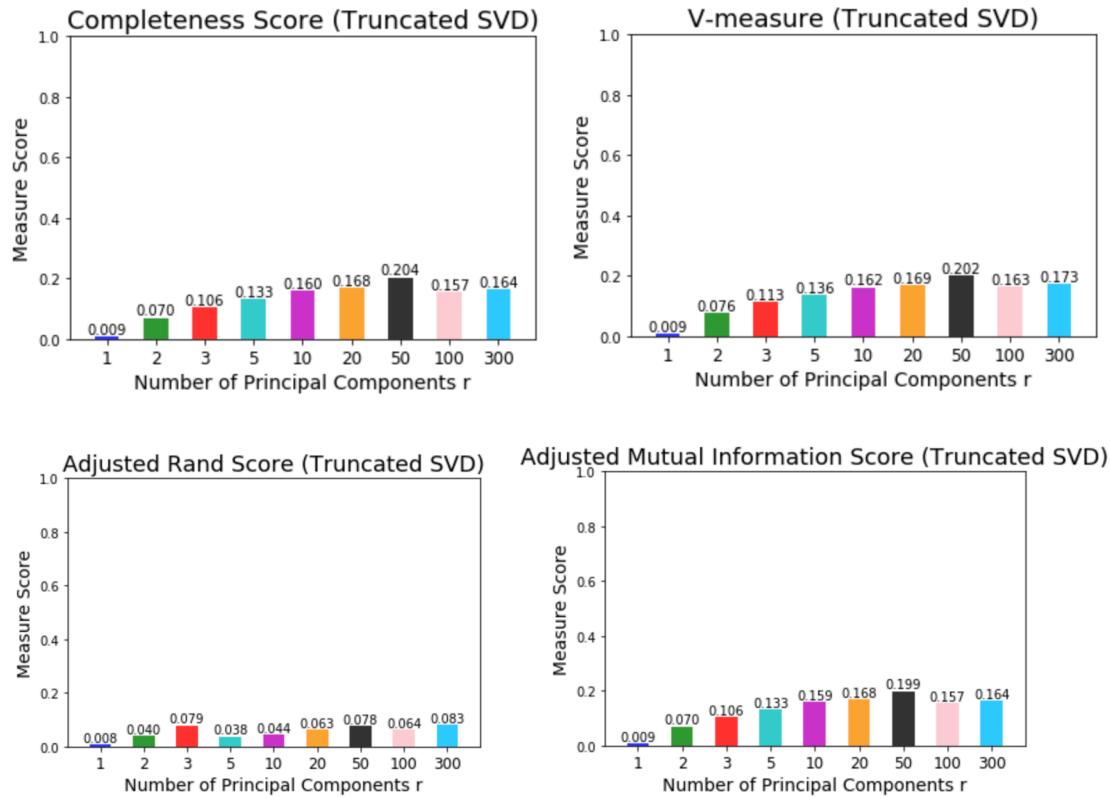
Class 1	Class 2	Class 3
comp.graphics	rec.autos	sci.crypt
comp.os.ms-windows.misc	rec.motorcycles	sci.electronics
comp.sys.ibm.pc.hardware	rec.sport.baseball	sci.med
comp.sys.mac.hardware	rec.sport.hockey	sci.space
Comp.windows.x		
Class 4	Class 5	Class 6
misc.forsale	talk.politics.misc	talk.religion.misc
	talk.politics.guns	alt.atheism
	talk.politics.mideast	soc.religion.christian

Below are our results:

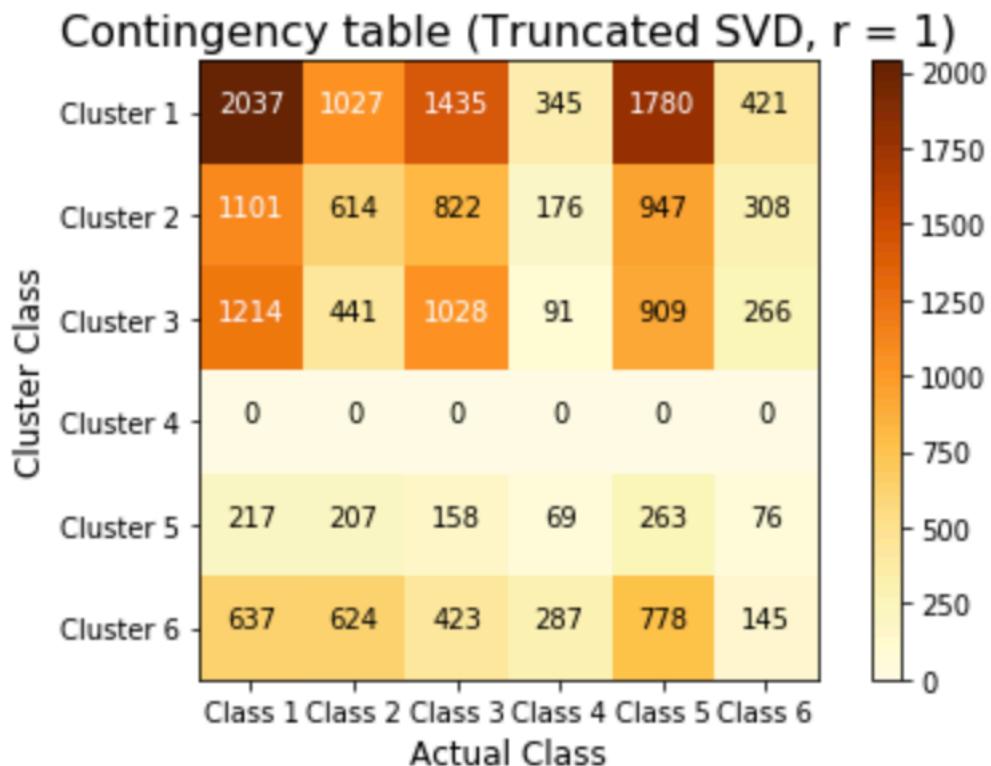
1. Truncated SVD + Normalizing

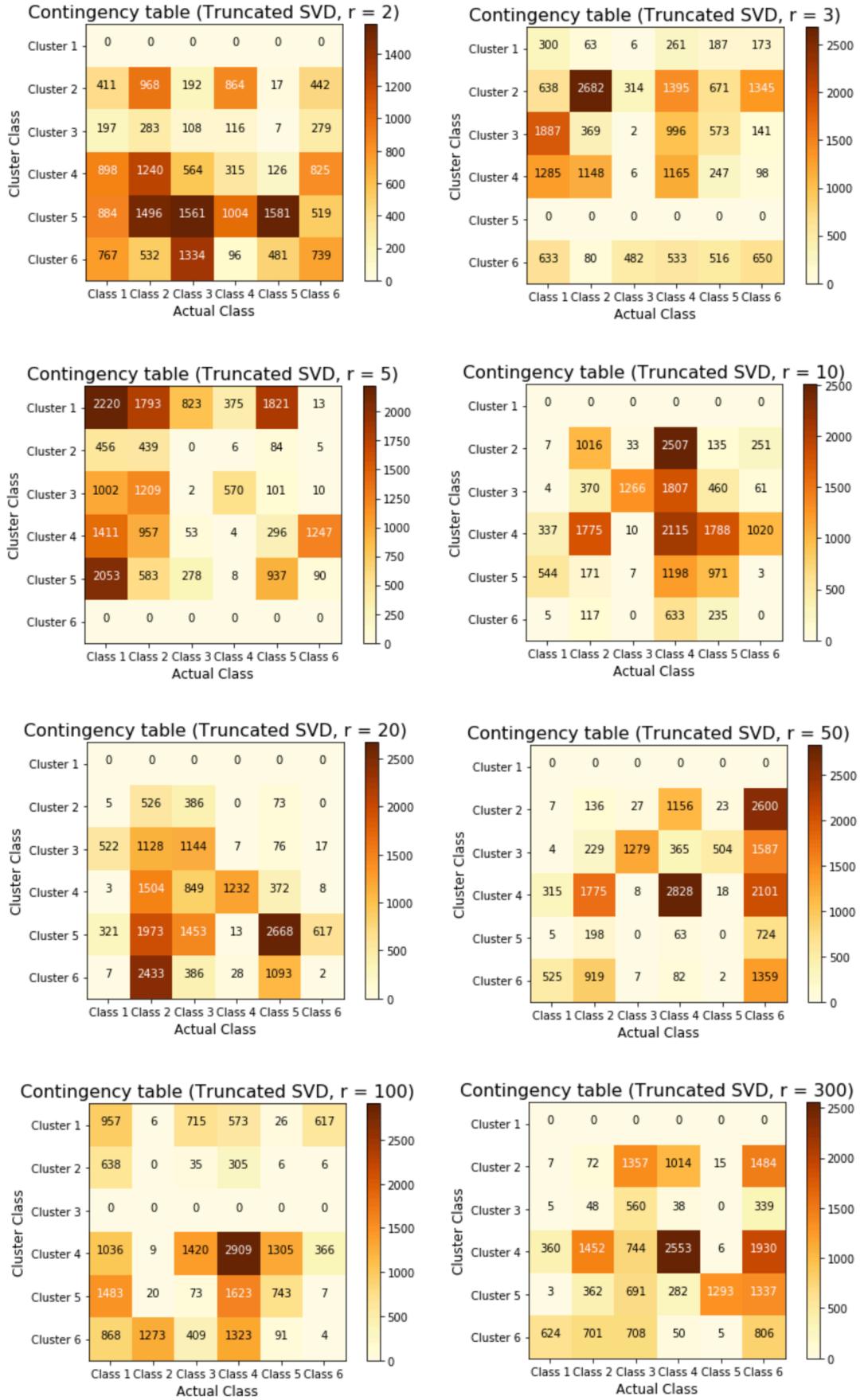
1) 5 measures scores



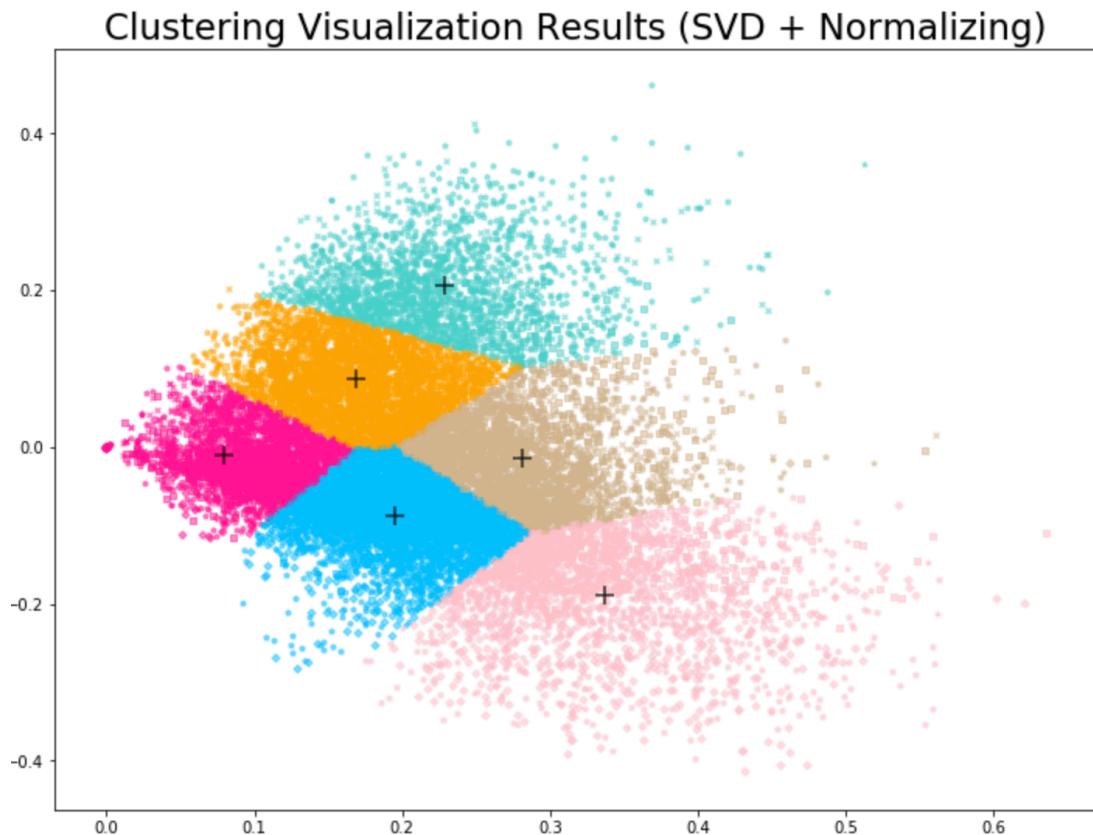


2) Contingency tables



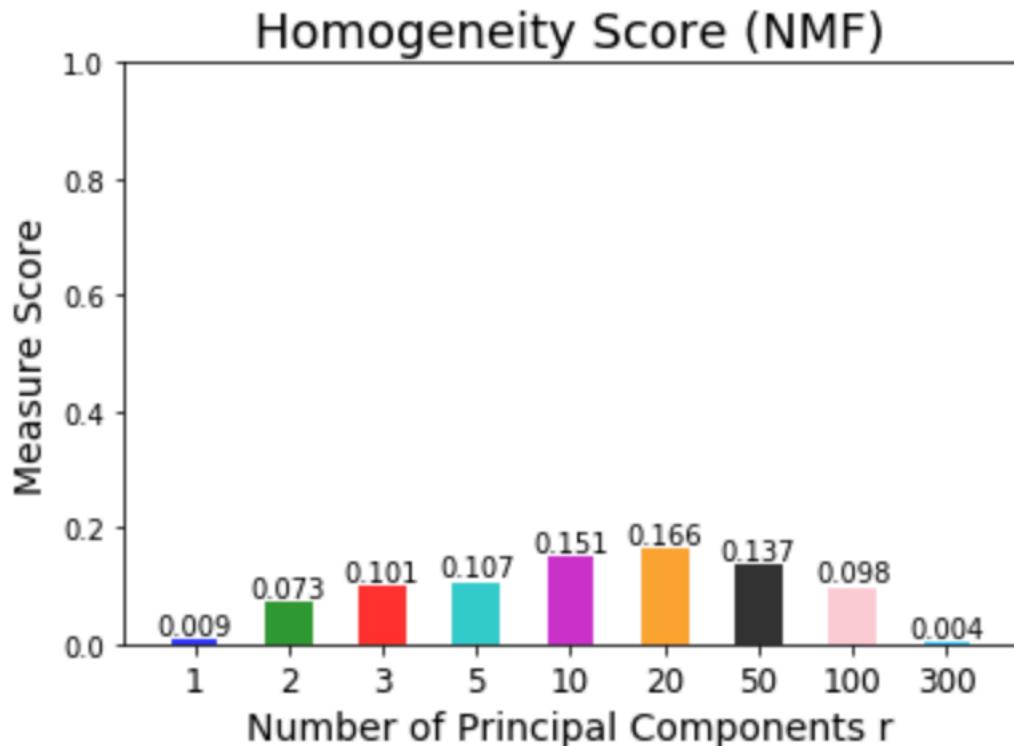


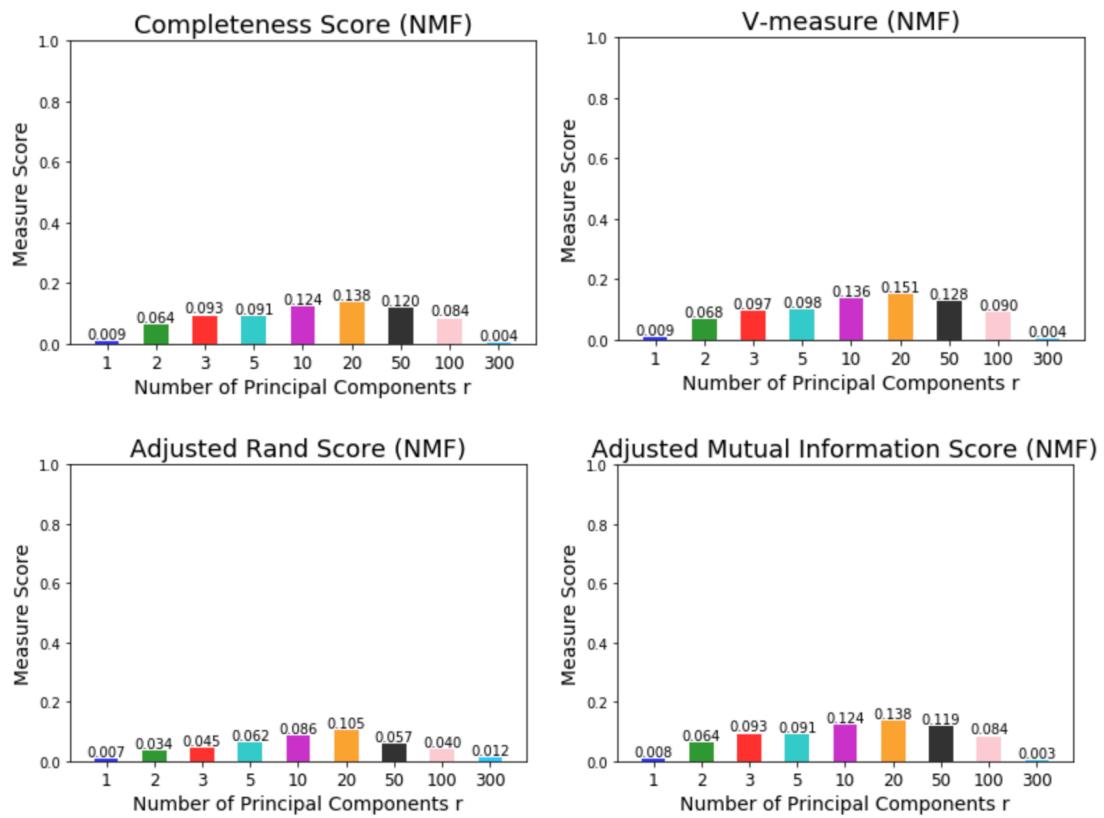
3) Clustering visualization results



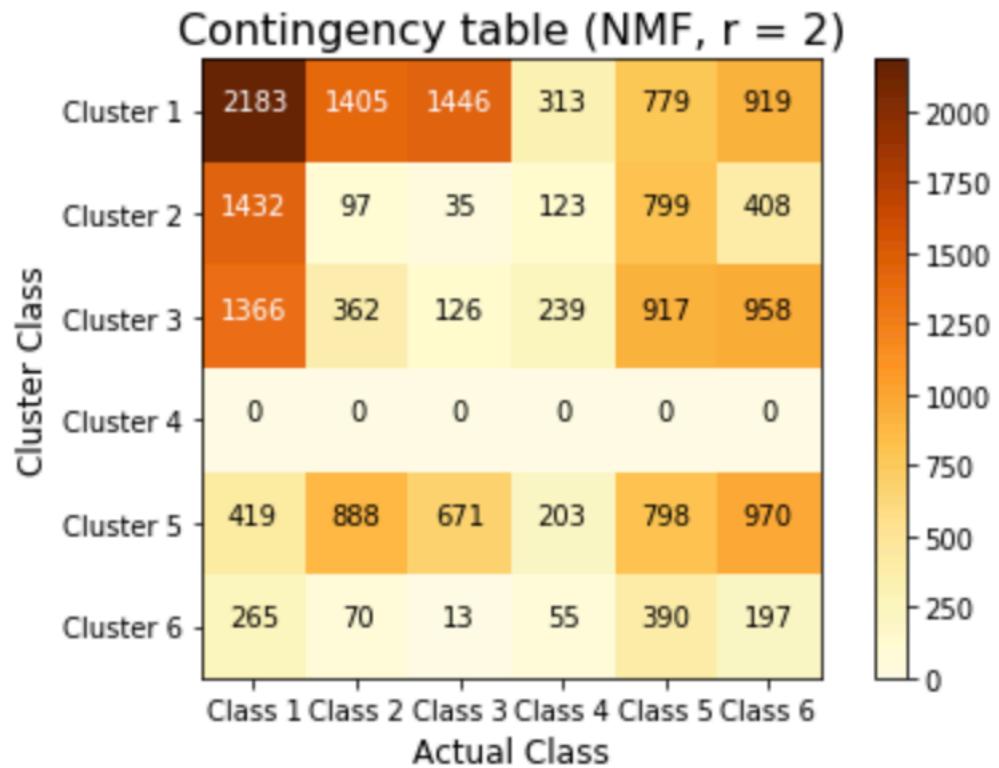
2. NMF + non-linear transformation

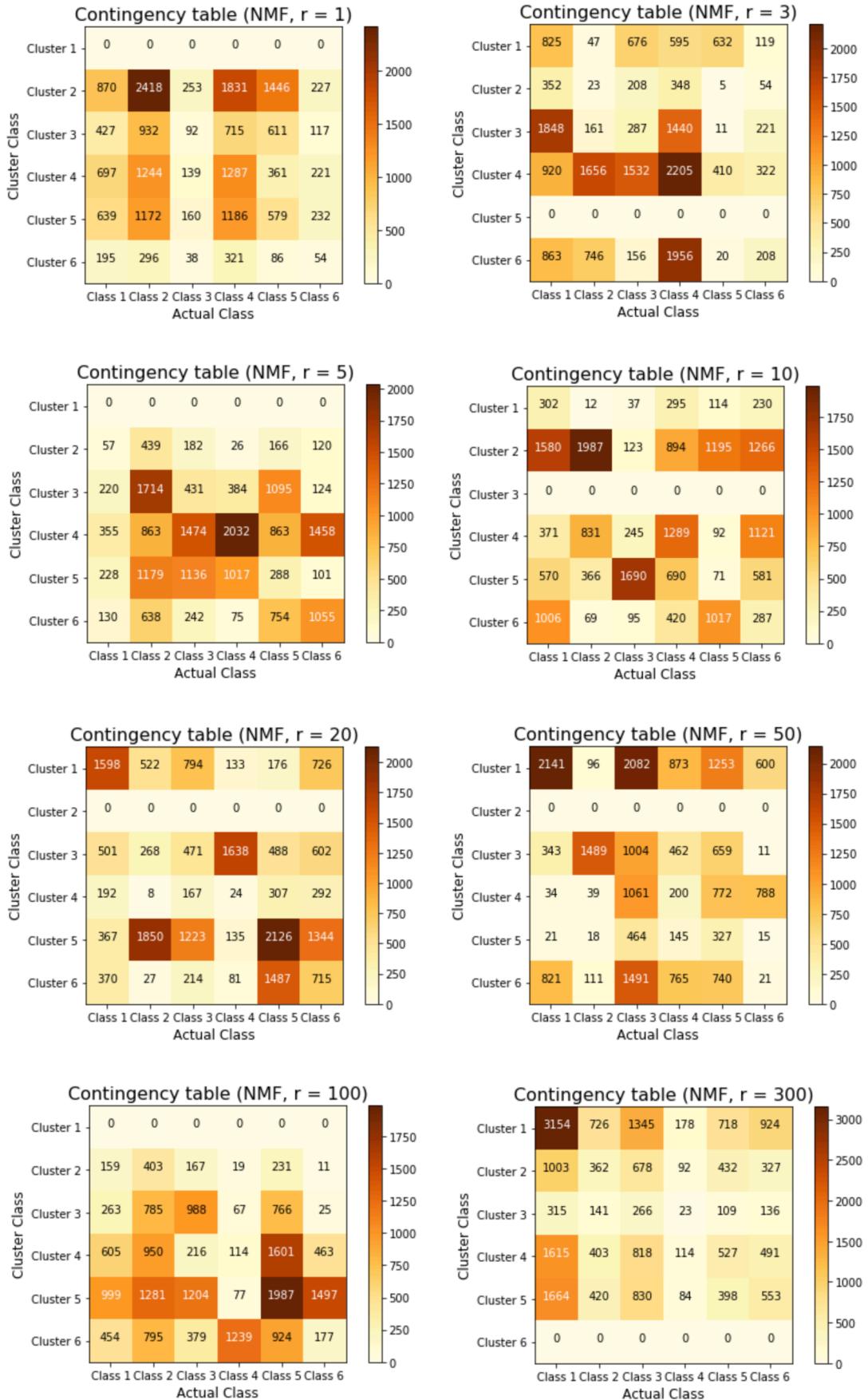
1) 5 measures scores



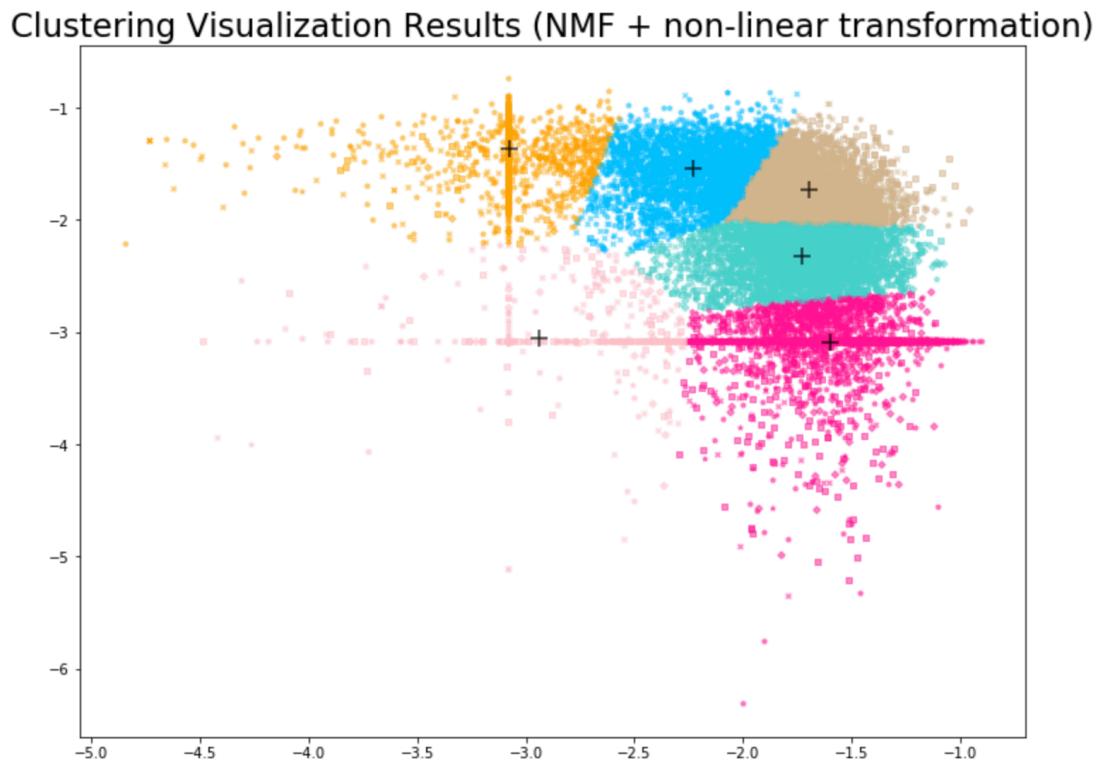


2) Contingency tables





3) Clustering visualization results



We can see from the results above, that the optimized reduced dimensions for the two methods we tried were different from previous, the optimal dimension for truncated SVD is 50 and 20 for NMF. And there exists an all-zero rows in each contingency table because that class 4 only has one subclass and the number of documents may too small compared to other classes to cluster as a new class.