# ECE 219 Project2

Jiuru Shao
UID:204288539
Haoxiang Zhang
UID:104278461

Feb 12th 2018

## 1 Introduction

In this experiment, we implemented a simple text classifier with unsupervised method(K-means clustering). Clustering differs from classification in that no *a priori* labeling (grouping) of the data points is available.

We represent of text file as TF-IDF matrix, and perform SVD and NMF on the matrix to reduce dimension. We evaluates our result with five metrics, namely the homogeneity score, the completeness score, the V-measure, the adjusted Rand score and the adjusted mutual info score.

We tried our code on both binary classification and multi-label classification. We find that the binary clustering works fine, but the 20-class clustering result are below our expectation despite all the methods we tried.

# 2 Problem 1. Building the TF-IDF Matrix

We use the same method as project 1 to perform TF-IDF, and the dimension of the TF-IDF matrix we get is $7882 \times 27768$.

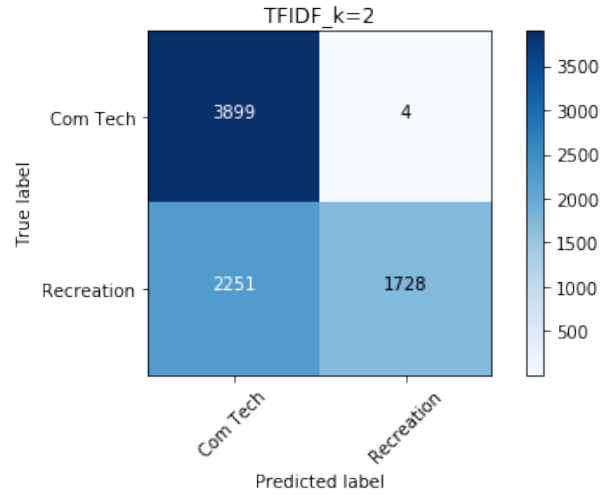# 3 Problem 2. Apply K-means Clustering using the TF-IDF

The results we get are:



Figure 1: **Contingency Matrix of TFIDF Data**

| | |
|---|---|
| Homogeneity | 0.255 |
| Completeness | 0.336 |
| V-measure | 0.290 |
| Adjusted rand score | 0.183 |
| Adjusted mutual info score | 0.255 |

Table 1: **5 Measures of TFIDF**

As we can see from figure 1 and table 1, the original TFIDF data provides a result that is not that bad. However, there are many cluster 2 data points are predicted as cluster 1 points, which is the reason that 5 measure scores are not very high.

# 4 Problem 3. Dimensionality Reduction

## 4.1 3a Report the percent of variance

In this part, we use the explained_variance_ratio in the TruncateSVD class to plot the graph. The result is in Figure 1.
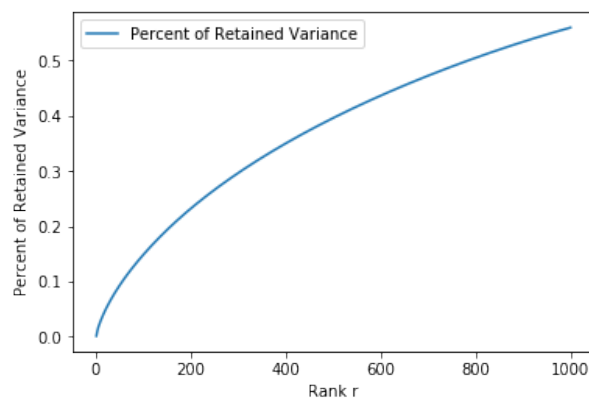


Figure 2: **The plot of the percent of variance the top $r$ principle components can retain v.s. $r$, for $r = 1$ to 1000.**

As $r$ keeps increasing, percent of retained variance increases as well.
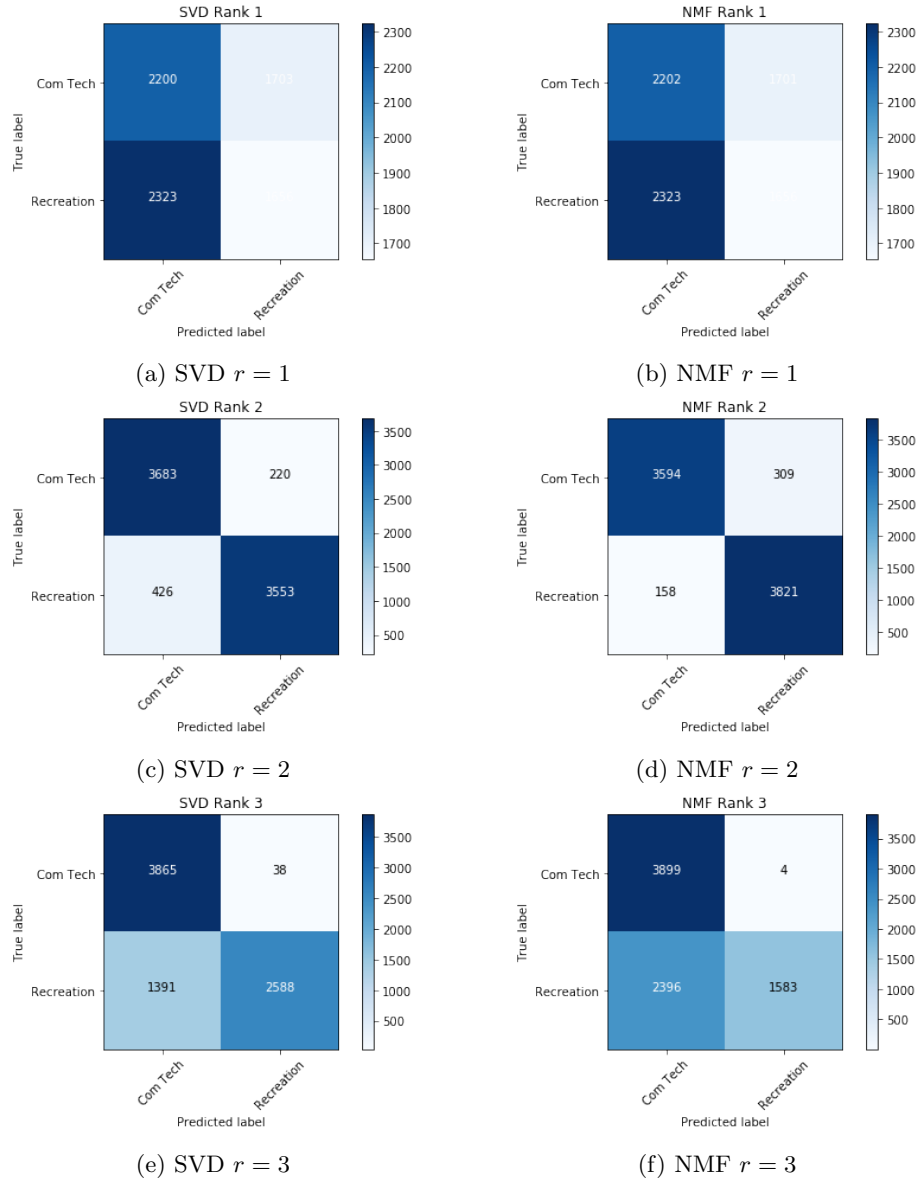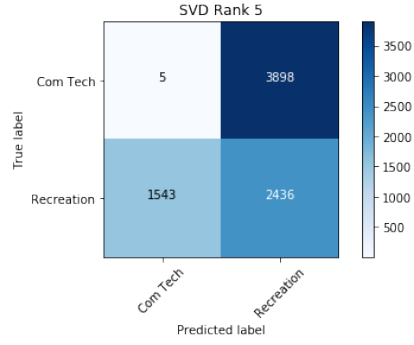
## 4.2   3b Apply SVD and NMF



(a) SVD $r = 1$

(b) NMF $r = 1$

(c) SVD $r = 2$

(d) NMF $r = 2$
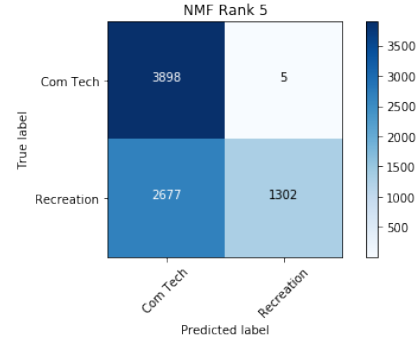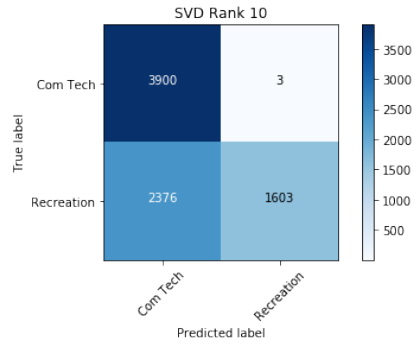
(e) SVD $r = 3$

(f) NMF $r = 3$

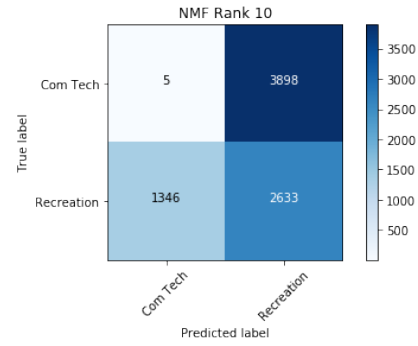Figure 3: **Contingency Matrix for** $r = 1, 2, 3$
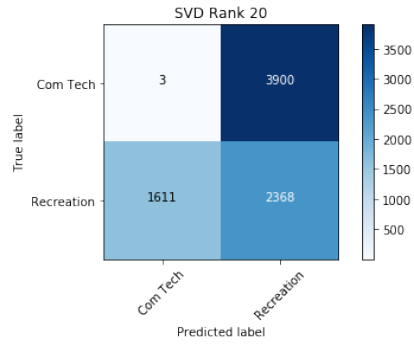
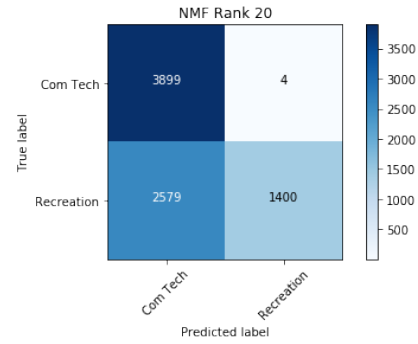(a) SVD $r = 5$            (b) NMF $r = 5$
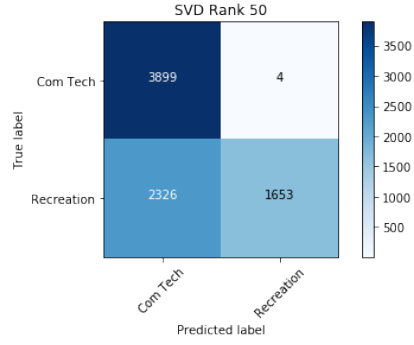
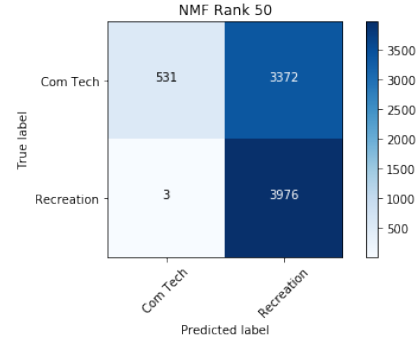(c) SVD $r = 10$          (d) NMF $r = 10$
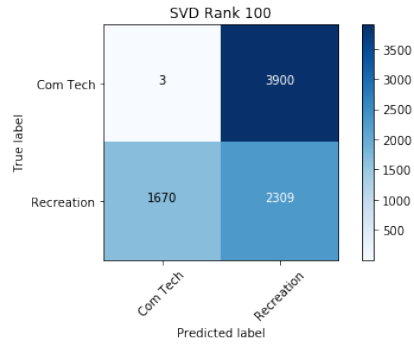
(e) SVD $r = 20$          (f) NMF $r = 20$

Figure 4: **Contingency Matrix for** $r = 5, 10, 20$
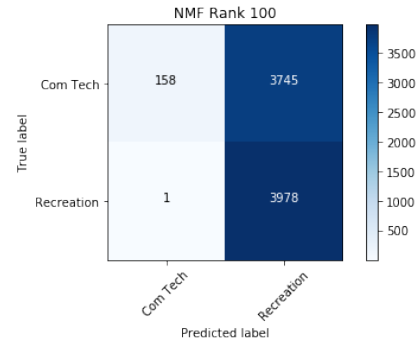
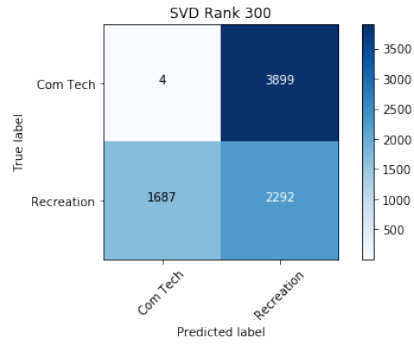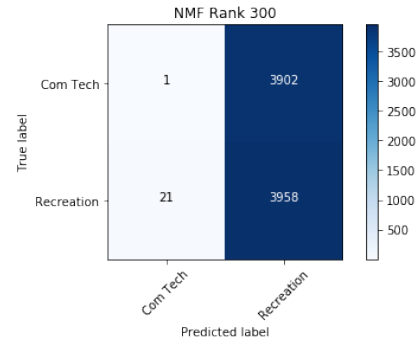(a) SVD $r = 50$        (b) NMF $r = 50$

(c) SVD $r = 100$        (d) NMF $r = 100$

(e) SVD $r = 300$        (f) NMF $r = 300$

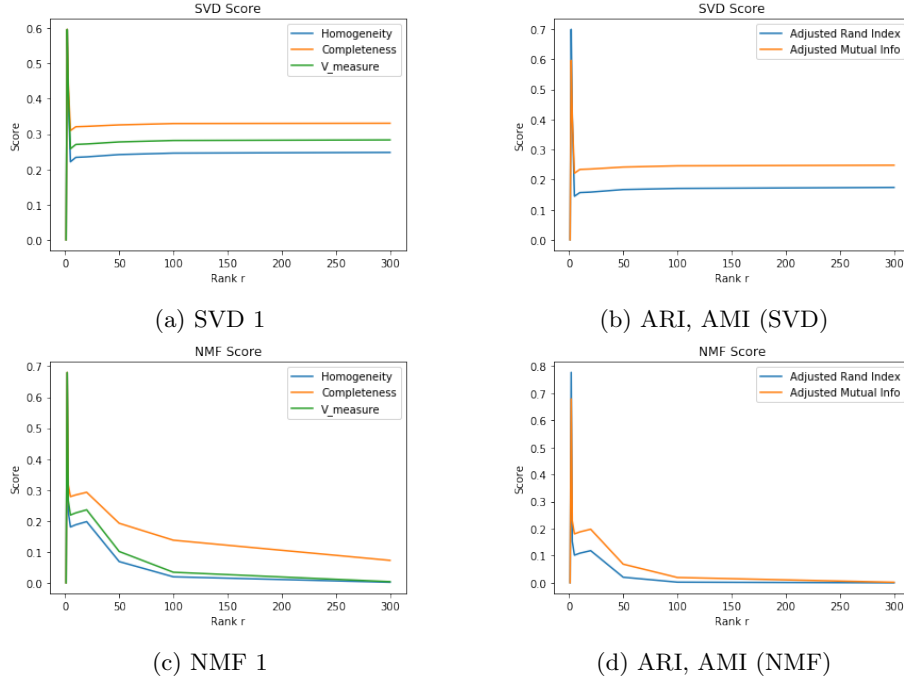Figure 5: **Contingency Matrix for** $r = 50, 100, 300$

(a) SVD 1

(b) ARI, AMI (SVD)

(c) NMF 1

(d) ARI, AMI (NMF)

Figure 6: **5 Measures**

From plots of contingency matrix and plots of 5 measures v.s. r, we conclude that the best $r$ choice for both SVD and NMF is $r = 2$.

We observe the non-monotonic behavior of the measure and one possible explanation is that Euclidean distance is not a good distance in high dimensions. In high dimensions, most of the mass of a multivariate Gaussian distribution is not near the mean, but in an increasingly distant "shell" around it; and most of the volume of a high-dimensional orange is in the skin, not the pulp.

# 5 Problem 4

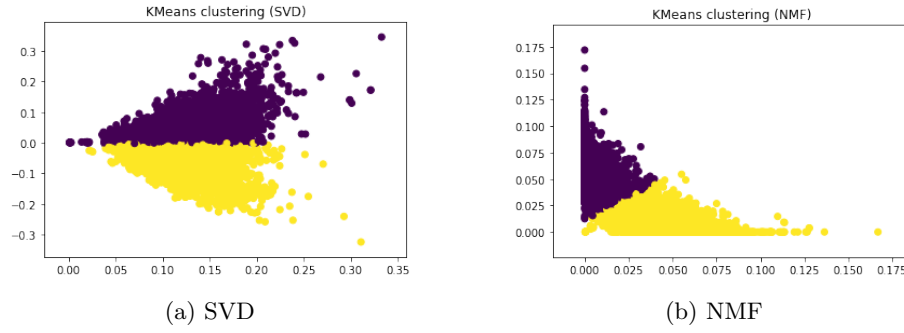## 5.1 4a Visualize Best Clustering Results



(a) SVD                                 (b) NMF

Figure 7: **Best Clustering Results**

|                            | SVD   | NMF   |
| -------------------------- | ----- | ----- |
| Homogeneity                | 0.596 | 0.679 |
| Completeness               | 0.597 | 0.680 |
| V-measure                  | 0.596 | 0.680 |
| Adjusted rand score        | 0.699 | 0.777 |
| Adjusted mutual info score | 0.596 | 0.679 |

Table 2: **5 Measure of Best Clustering Results**

From the distribution of the data points after SVD, we can observe that it is very hard to cluster them into 2 classes, because most of the points are very close to each other and there is no clear boundary. In contrast, NMF is slightly better, showing a triangular distribution. Table 2 confirms this observation.

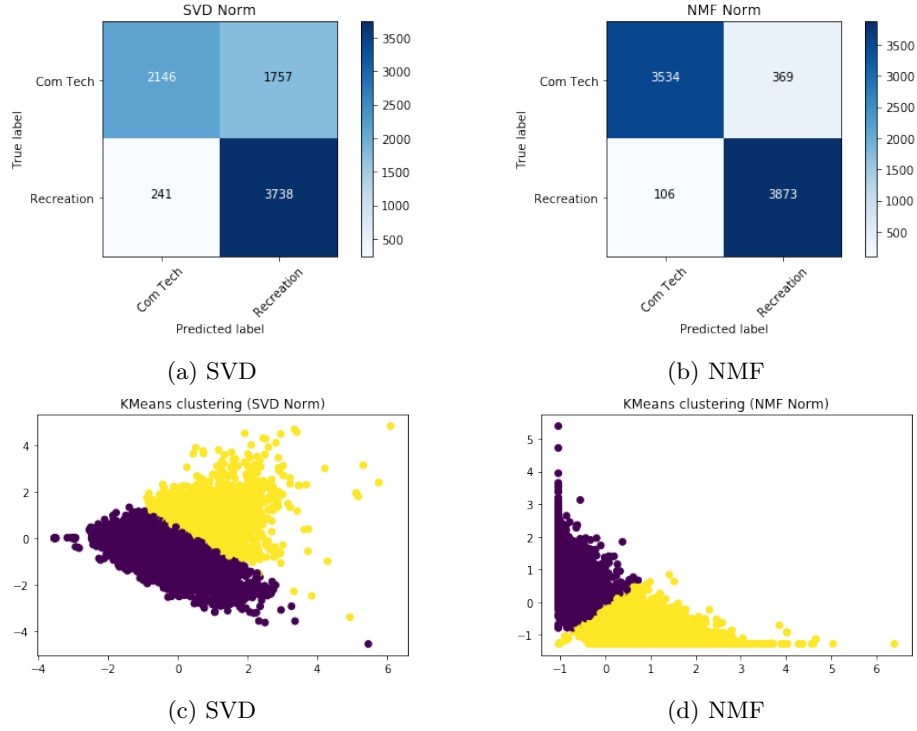## 5.2    4b Different Transformations

### 1. Normalization



(a) SVD

(b) NMF

(c) SVD

(d) NMF

Figure 8: **Plots after Normalization**

|  | SVD | NMF |
|---|---|---|
| Homogeneity | 0.235 | 0.683 |
| Completeness | 0.264 | 0.686 |
| V-measure | 0.249 | 0.684 |
| Adjusted rand score | 0.255 | 0.773 |
| Adjusted mutual info score | 0.235 | 0.683 |

Table 3: **5 Measures after Normalization**

We can see the result of normalization after SVD decomposition becomes significantly worse, while the result after NMF remains the same, we can confirm this by visualizing the data points: the distribution of the observation points still remain the similar shape.
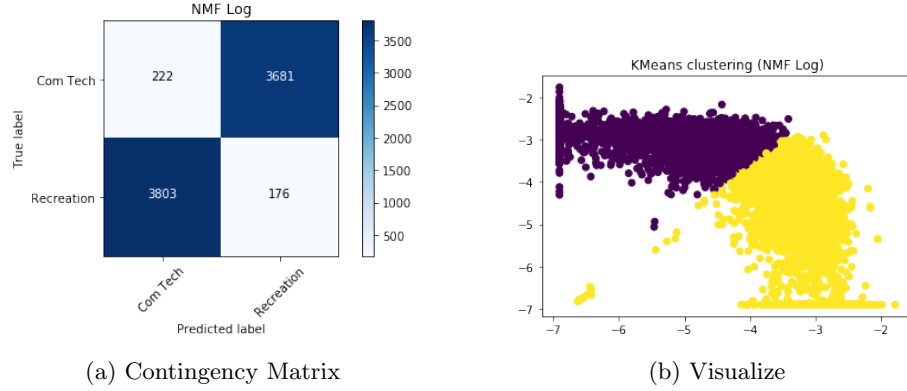
2. **Logarithm Transformation (NMF Only)**



(a) Contingency Matrix        (b) Visualize

Figure 9: **Plots after Log**

| | |
|---|---|
| Homogeneity | 0.712 |
| Completeness | 0.712 |
| V-measure | 0.712 |
| Adjusted rand score | 0.808 |
| Adjusted mutual info score | 0.712 |

Table 4: **Measures after Log**

We find log transformation greatly improves the result because the original data points are skewed with wide distribution, and taking the log of the features may restore symmetry to it. As we can see from the distribution of the observations after taking log, the points are mapped to a sector that is easier to cluster.
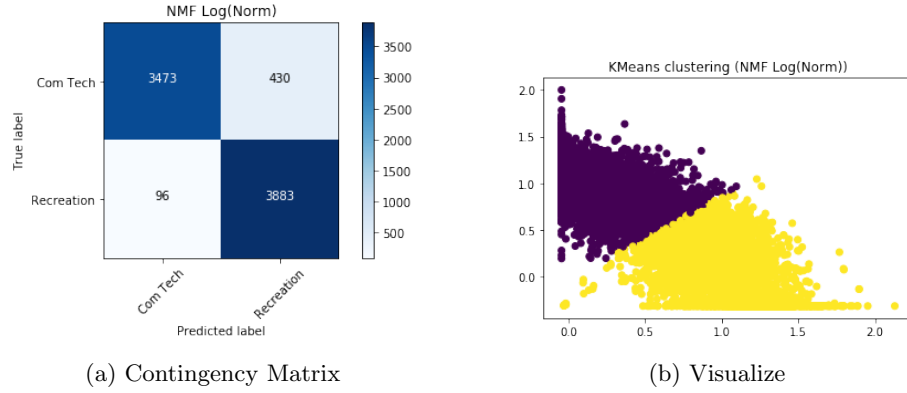
3. **Norm then Log (NMF Only)**



(a) Contingency Matrix

(b) Visualize

Figure 10: **Plots (Norm then Log)**

| | |
|---|---|
| Homogeneity | 0.663 |
| Completeness | 0.667 |
| V-measure | 0.665 |
| Adjusted rand score | 0.751 |
| Adjusted mutual info score | 0.663 |

Table 5: **5 Measures (Norm then Log)**

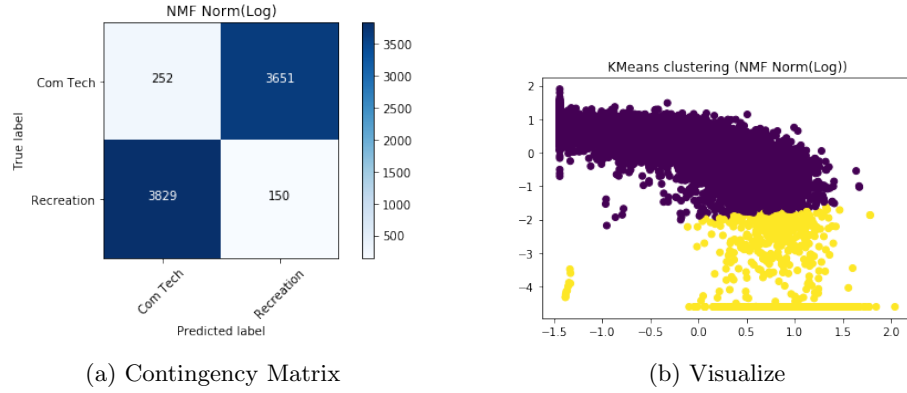The result is roughly the same as part 4a, there is no improvement.

4. **Log then Norm (NMF Only)**



(a) Contingency Matrix          (b) Visualize

Figure 11: **Plots (Log then Norm)**

| Homogeneity | 0.711 |
|---|---|
| Completeness | 0.712 |
| V-measure | 0.712 |
| Adjusted rand score | 0.806 |
| Adjusted mutual info score | 0.711 |

Table 6: **5 Measures (Log then Norm)**

The result is similar to taking log only.

# 6 Problem 5

## 6.1 Get TF-IDF matrix

The dimension of the TF-IDF matrix is $18846 \times 52295$.

## 6.2 Apply K-means clustering with k = 20



Figure 12: **Contingency Matrix of TFIDF Data**

| Homogeneity | 0.322 |
|---|---|
| Completeness | 0.398 |
| V-measure | 0.356 |
| Adjusted rand score | 0.119 |
| Adjusted mutual info score | 0.320 |

Table 7: **5 Measures of TFIDF**

As we can see from figure 12 and table 7, the original TFIDF data provides a result that is not bad because this is a clustering for 20 different clusters.

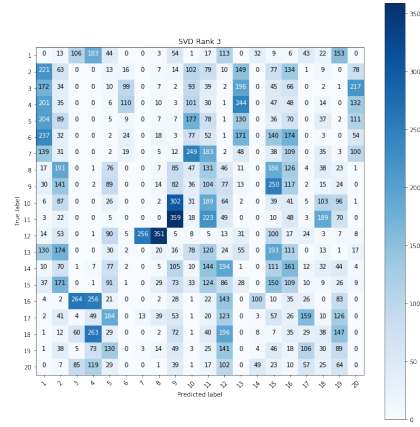## 6.3 Dimension reduction with SVD and NMF

(a) SVD $r = 1$

(b) NMF $r = 1$

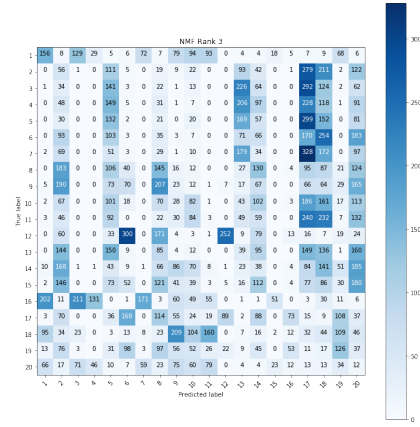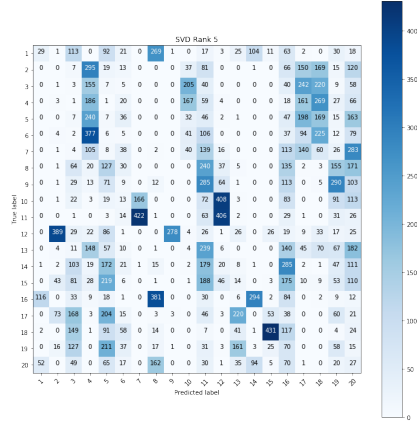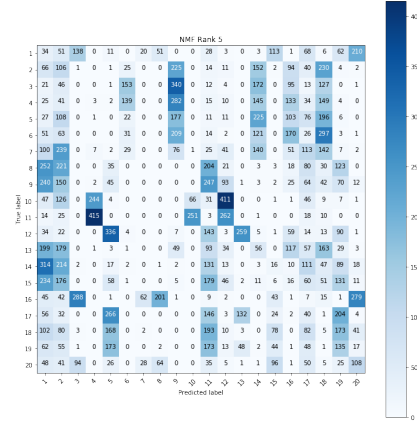(c) SVD $r = 2$
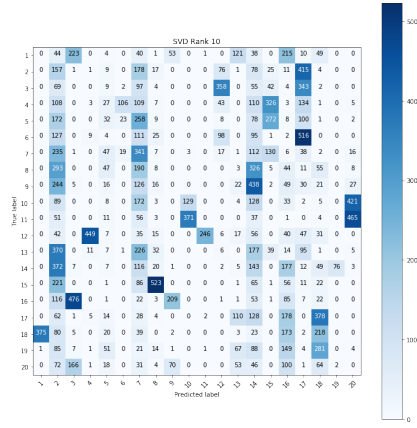
(d) NMF $r = 2$

(e) SVD $r = 3$

(f) NMF $r = 3$

Figure 13: **Contingency Matrix for** $r = 1, 2, 3$
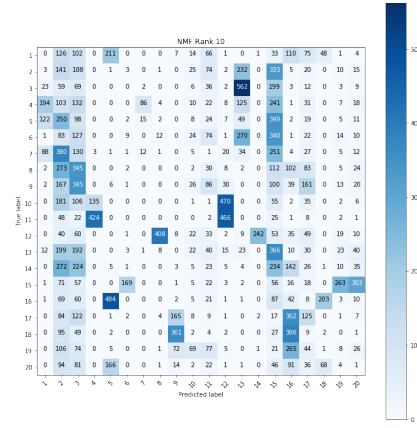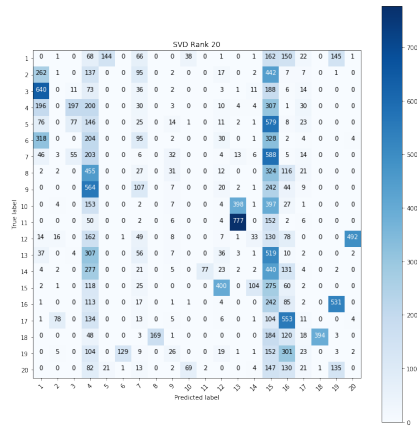
14

(a) SVD $r = 5$          (b) NMF $r = 5$
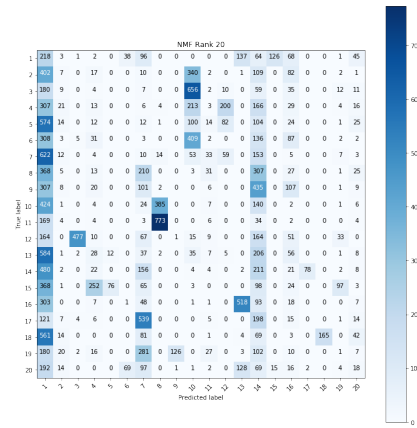
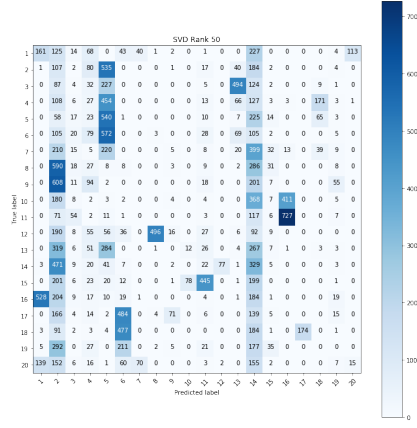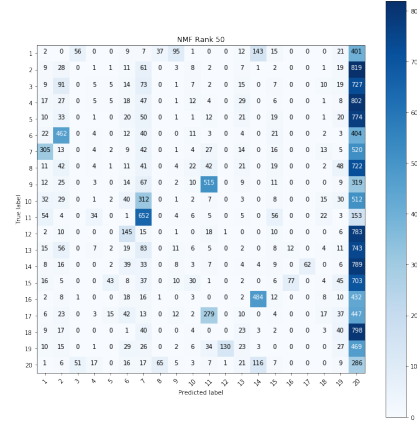(c) SVD $r = 10$          (d) NMF $r = 10$

(e) SVD $r = 20$          (f) NMF $r = 20$

Figure 14: **Contingency Matrix for** $r = 5, 10, 20$

15

(a) SVD $r = 50$

(b) NMF $r = 50$

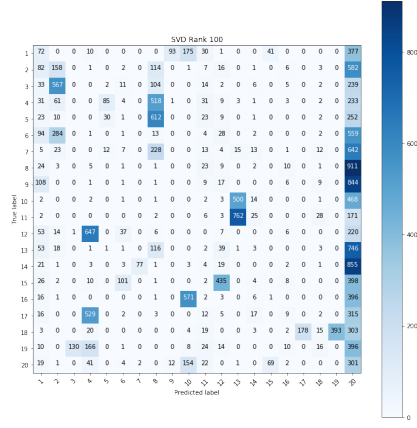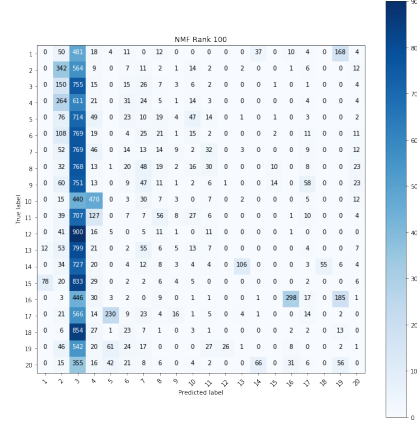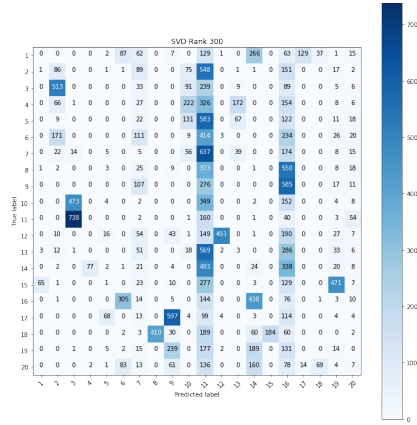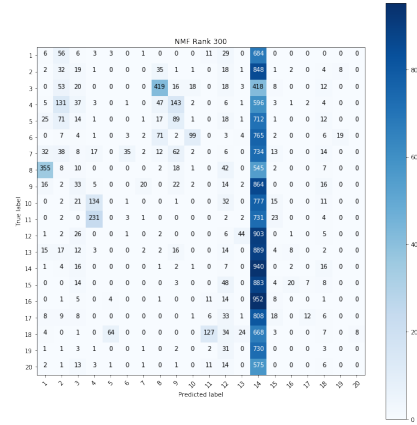(c) SVD $r = 100$

(d) NMF $r = 100$
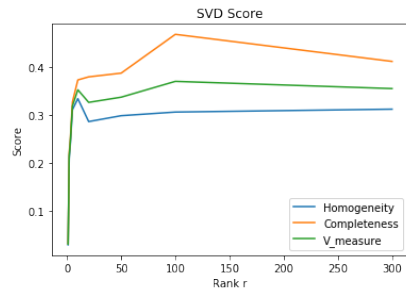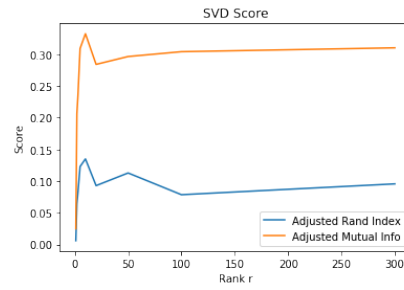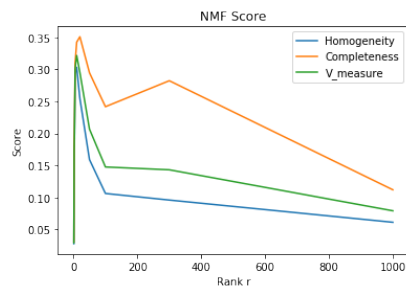
(e) SVD $r = 300$

(f) NMF $r = 300$

Figure 15: **Contingency Matrix for** $r = 50, 100, 300$

(a) SVD 1

(b) ARI, AMI (SVD)

(c) NMF 1

(d) ARI, AMI (NMF)

Figure 16: **5 Measures of 20 Clusters**

17

As in problem 3, we explored the performance measures for $r = 1, 2, 3, 5, 10,$ $20, 50, 100, 300$. For NMF, we expaned the performance measures to $r = 1000$ as well to observe the trend. However we find that unlike the previous part, the optimal $r$ for each metric are not the same at each time. We pick adjusted mutual info score as the metric, and the result shows $r = 10$ is the best value for both SVD and NMF. The details of the measures at $r = 10$ are:

|  | SVD | NMF |
| --- | --- | --- |
| Homogeneity | 0.335 | 0.317 |
| Completeness | 0.374 | 0.357 |
| V-measure | 0.353 | 0.336 |
| Adjusted rand score | 0.135 | 0.124 |
| Adjusted mutual info score | 0.332 | 0.315 |

Table 8: **Measures of Best Clustering Results**

We visualize the best clustering results of 20 clusters:



(a) SVD

(b) NMF

Figure 17: **Best Clustering Results**

18

## 6.4    Apply Normalization and Log

### 1.  Normalization
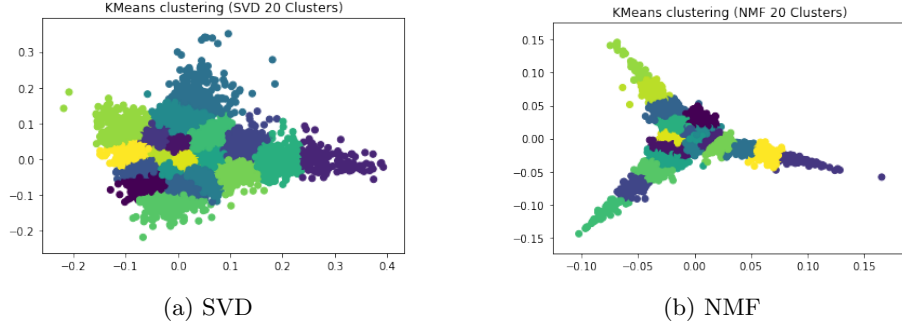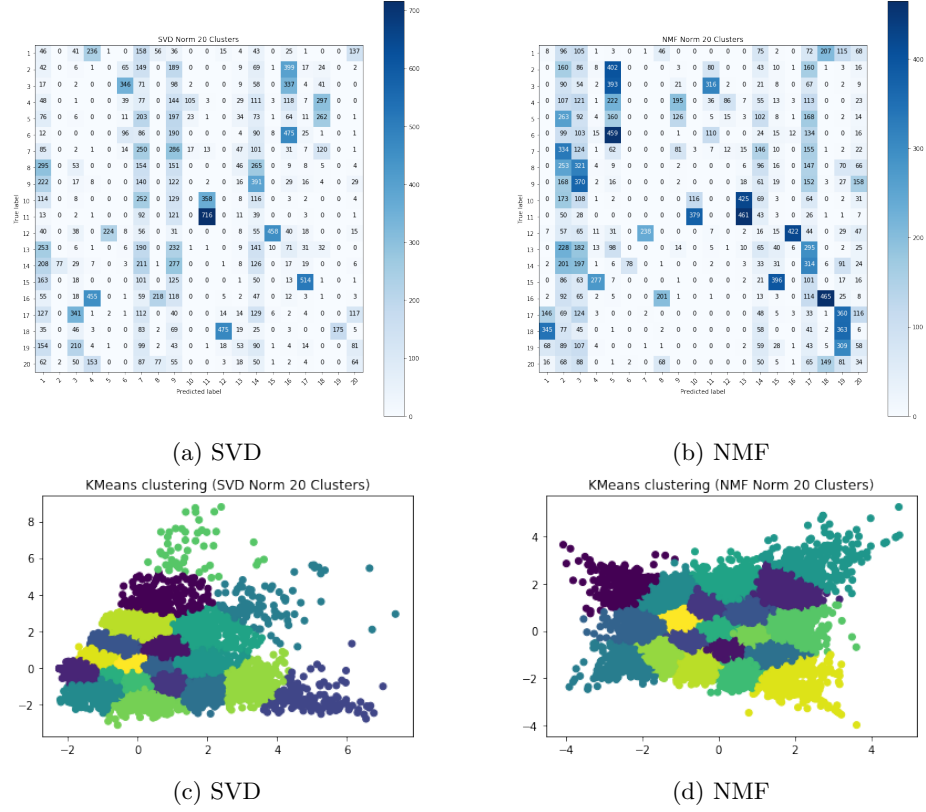


(a) SVD

(b) NMF



(c) SVD

(d) NMF

Figure 18: **Plots after Normalization**

|  | SVD | NMF |
|---|---|---|
| Homogeneity | 0.310 | 0.305 |
| Completeness | 0.349 | 0.341 |
| V-measure | 0.328 | 0.322 |
| Adjusted rand score | 0.127 | 0.117 |
| Adjusted mutual info score | 0.308 | 0.303 |

Table 9: **5 Measures after Normalization**

Despite the changes in the pattern(the data points become more "spread-out") of distribution after normalization, the result is still not better than before.

2. **Log Transformation with NMF**



(a) Contingency Matrix   (b) Visualize

Figure 19: **Plots after Log**

| Homogeneity | 0.376 |
|---|---|
| Completeness | 0.380 |
| V-measure | 0.378 |
| Adjusted rand score | 0.205 |
| Adjusted mutual info score | 0.374 |

Table 10: **Measures after Log**

### 3. Norm then Log with NMF



(a) Contingency Matrix

(b) Visualize

Figure 20: **Plots (Norm then Log)**

| Homogeneity | 0.337 |
|---|---|
| Completeness | 0.348 |
| V-measure | 0.342 |
| Adjusted rand score | 0.157 |
| Adjusted mutual info score | 0.335 |

Table 11: **5 Measures (Norm then Log)**

4. **Log then Norm with NMF**



(a) Contingency Matrix

(b) Visualize
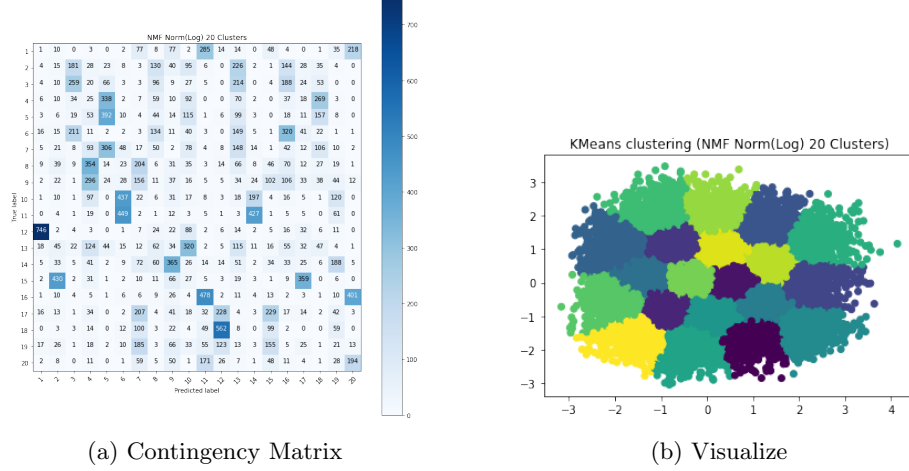
Figure 21: **Plots (Log then Norm)**

| Homogeneity | 0.368 |
| --- | --- |
| Completeness | 0.370 |
| V-measure | 0.369 |
| Adjusted rand score | 0.202 |
| Adjusted mutual info score | 0.366 |

Table 12: **5 Measures (Log then Norm)**

From our experiment, we can see that the best result is generated with the log transformation only. We also notice that even the best result only have a minor improvement on the metrics despite the transformation in the distribution of data points.

# 7 Conclusion

In summary, we achieve the goal of this project. We used TF-IDF and SVD/NMF to effectively represent the features and reduce the dimensions. In the case of binary classification, our result shows that NMF is more effective than SVD if we perform a log transformation after dimension reduction. Also, a log transformation followed by normalization will yield similar result.

On the other hand, we also tried to expand our dataset by clustering the 20 original data sets. However, the result is much worse than the binary case. Similarly, the best result is achieved with log transformation with NMF. This

indicates that our method is not able to cluster the texts into more refined classes effectively.