# ECE 219 Project 2

Tanner Dulay, Loïc Maxwell, Zoë Tucker

February 3, 2021

## 1   Question 1

In this assignment, we explore various methods for clustering of text data. To begin, we use the "20 Newsgroups" dataset, taking the two well-separated classes previously defined ("computer technology" and "recreational activity"). As this project uses unsupervised learning, we do not distinguish between training and test data, and instead use the entire dataset throughout the project.

To start our process, we build the TF-IDF matrix. We use `min_df` = 3, exclude English stopwords/punctuation, and remove headers and footers, but we do not perform stemming or lemmatization. The resulting TF-IDF matrix is **7882 × 21909**.

## 2   Question 2

Next, we apply k-means clustering directly to the TF-IDF data. As we are interested in two classes of documents, we set $k = 2$. The below contingency matrix (Figure 1) shows the results of this.
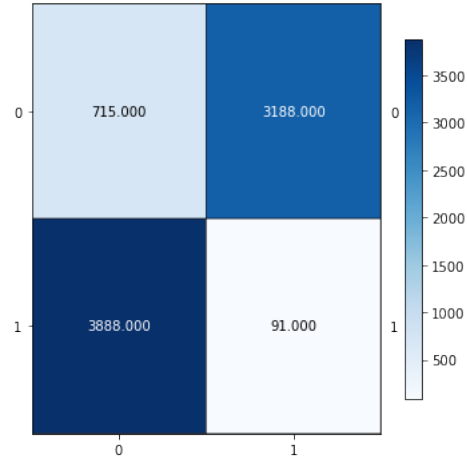
Figure 1: Contingency Matrix for K-Means on TF-IDF Matrix

# 3 Question 3

We then find the metrics we need to evaluate our classification results. We find the following values:

| | |
|---|---|
| Homogeneity score | 0.5600025412974119 |
| Completeness score | 0.5716557539454379 |
| V_measure score | 0.5657691483417876 |
| Adjusted Rand Index score | 0.6327475635535726 |
| Adjusted Mutual Information Score | 0.5657289846874574 |

# 4 Question 4

As the metrics found in Question 3 do not demonstrate a good clustering result, we now try dimensionality reduction to help our data better fit the assumptions made by k-means algorithms.

We begin by finding the percent of variance the top $r$ principal components can retain for various values of $r$. Figure 2 shows our findings.
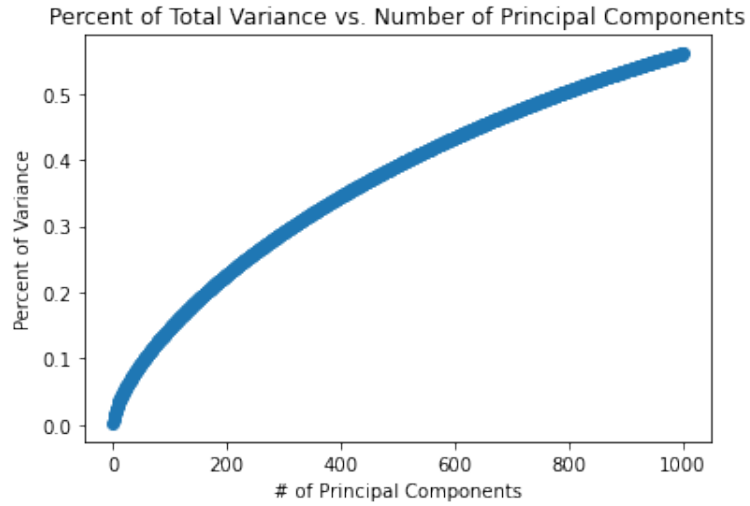
Figure 2: Principal Component Results

# 5  Question 5

Next, we will try to experimentally find the value of $r$ (i.e. the number of dimensions) that we can reduce our data to in order to obtain the best clustering results. Testing $r \in \{1, 2, 3, 5, 10, 20, 50, 100, 300\}$ using both NMF and SVD, we plot the values of the five metrics found in Question 3 for each $r$. The results are shown in Figures 3 and 4.
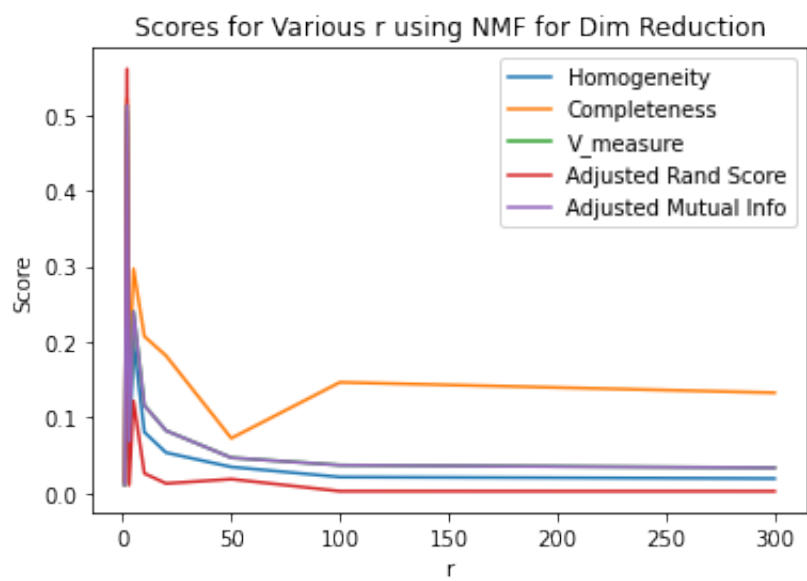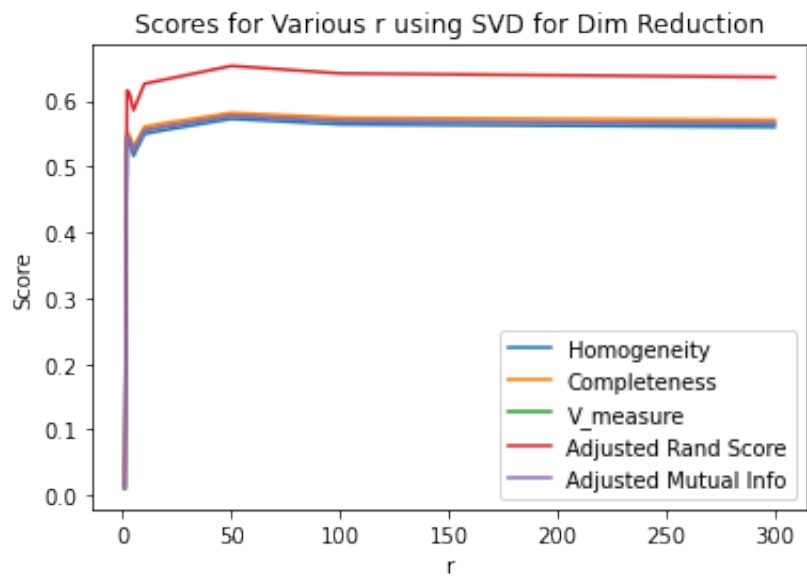
Figure 3: NMF Results



Figure 4: SVD Results

# 6    Question 6

As we can see, the behavior of the measures as $r$ increases is non-monotonic.

This is due to a trade-off between a larger value of $r$ giving a more faithful representation of the data, and a smaller value of $r$ being better suited for k-means clustering. The higher the dimensionality of a dataset, the more difficult it is to meaningfully evaluate the "distance" between data points, which is what k-means relies on. Thus, we expect there to be an optimal value of $r$ that balances these two concerns and maximizes our scores. In our case, we find that $r = 2$ for NMF, and $r = 50$ for SVD.

# 7    Question 7

Now, we want to visualize our clustering results. We do this by using k-means clustering with our dataset reduced by NMF and SVD with the optimal values of $r$ found in Question 6, then projecting the data into 2-dimensional space with SVD and coloring the points according to their label. Figures 5 and 6 show the results for NMF with $r = 2$. Figures 7 and 8 show the results for SVD with $r = 50$.
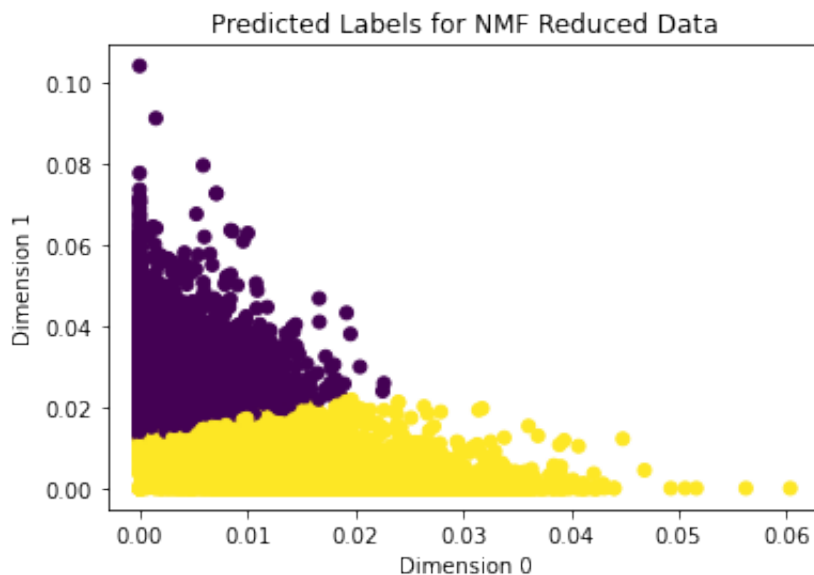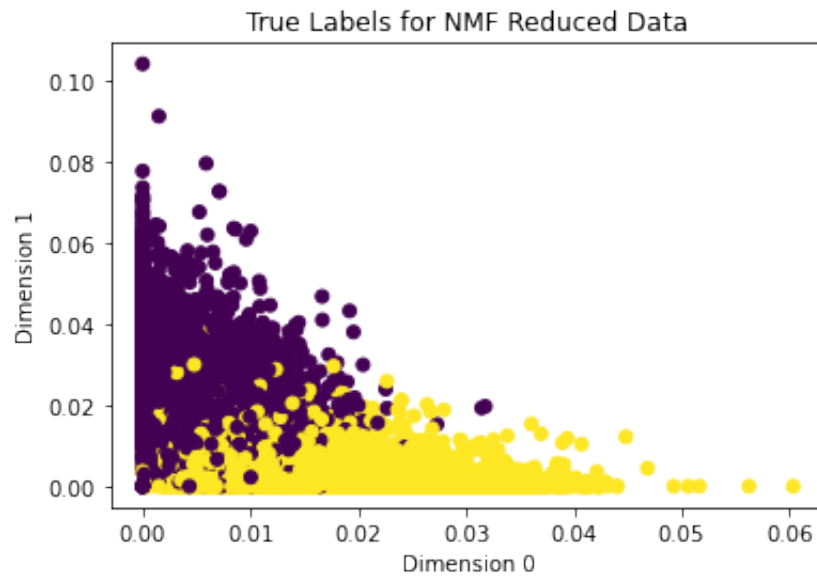


Figure 5: Predicted labels for NMF reduced data

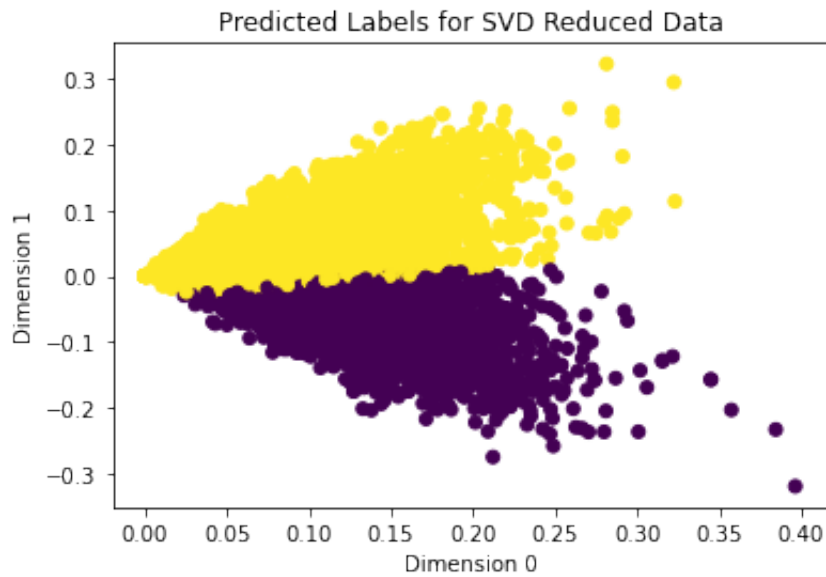Figure 6: Ground truth labels for NMF reduced data
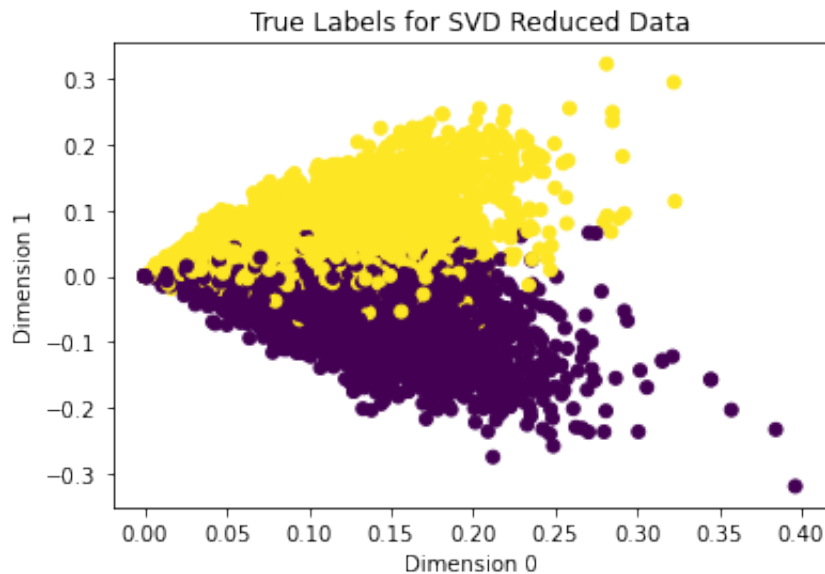


Figure 7: Predicted labels for SVD reduced data

Figure 8: Ground truth labels for SVD reduced data

# 8    Question 8

Because the NMF reduced data was already in 2D, we did not reduce it further with SVD. However, if we were to do SVD reduction, the resulting matrix would still be positive, as it is in Figures 5 and 6. SVD is an orthogonal projection that preserves the angles between feature vectors.

The 2D visualization from Figure 8 shows that the SVD reduced data looks somewhat linearly separable around a value of 0 for Dimension 1, which stems from the second principal component of the TF-IDF matrix. This suggests that this principal component holds important information pertaining to the class that the data belongs to. Although this analysis is restricted to two dimensions, this also suggests that this type of distribution would be good for KMeans clustering. A similar type of separation can be found at the diagonal between dimensions 0 and 1 for the NMF reduced data in Figure 6. This would also make it ideal for KMeans clusering.

However, for both types of reductions, because there is no way to perfectly linearly separate the data there are some better options out there for properly classifying this data. Perhaps a classifier that can learn non-linear behavior in the data, such as a neural network with ReLU/eLU, can perform better on this dataset.

# 9 Question 9

Using SVD with r=50, which gave us the best results on the two classes, we got the following contingency matrix for KMeans on the entire dataset of 20 classes:
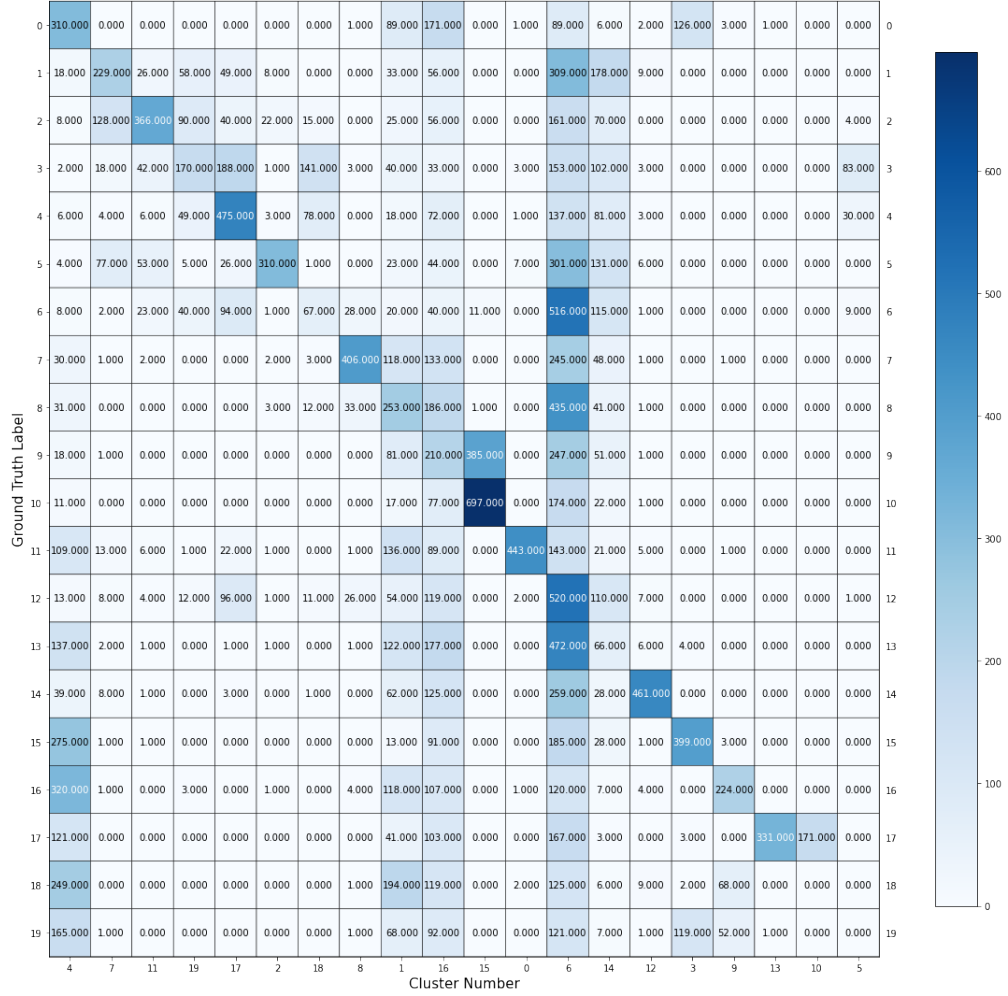


Figure 9: Contingency Matrix for KMeans on 20 Class Data

The respective scores for this clustering are as follows:

| | |
|---|---|
| Homogeneity score | 0.31630458145750995 |
| Completeness score | 0.3682552023183464 |
| V_measure score | 0.3403086491478565 |
| Adjusted Rand Index score | 0.10379799750286293 |
| Adjusted Mutual Information Score | 0.3380068278889245 |

# 10    Question 10

Using the Kullback-Leibler Divergence loss metric for NMF and r=10, we were able to get the following scores:

| | |
|---|---|
| Homogeneity score | 0.37714629232162333 |
| Completeness score | 0.3900505160910544 |
| V_measure score | 0.3834898799077068 |
| Adjusted Rand Index score | 0.18833101982676798 |
| Adjusted Mutual Information Score | 0.38146535004967785 |

There is indeed a noticeable improvement over the performance using ordinary SVD with r=50 (from Q9).

# 11    Question 11

We compared the performance of KMeans on the dataset reduced with UMAP using various values for n_components. The plots are shown below:
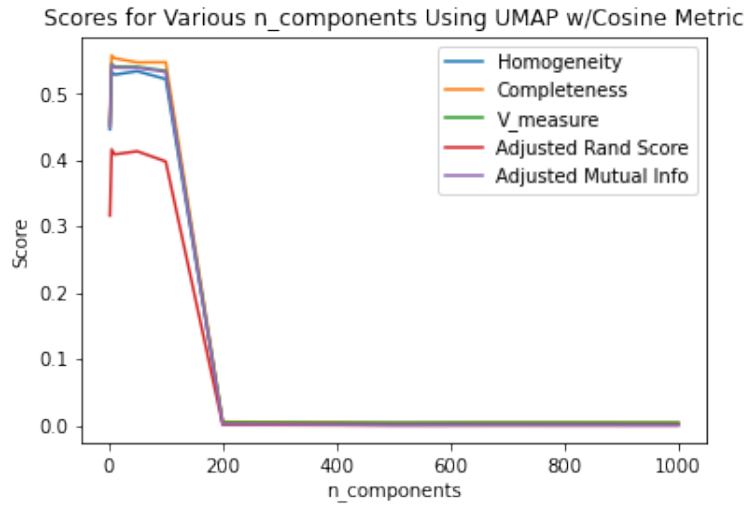


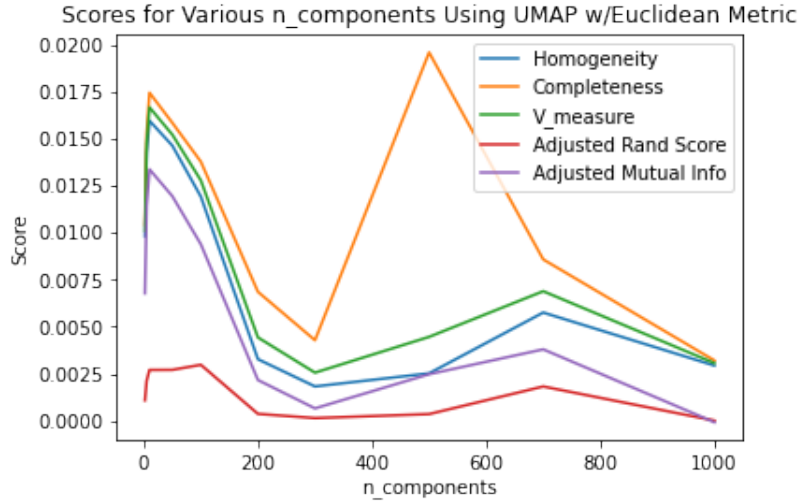Figure 10: Metric scores for Cosine Metric vs. n_components

Figure 11: Metric scores for Euclidean Metric vs. n_components

The best scores using both metrics came from **n_components=5**:

| Cosine Metric | |
|---|---|
| Homogeneity score | 0.5347926404764849 |
| Completeness score | 0.5549487009354954 |
| V_measure score | 0.5446842655666289 |
| Adjusted Rand Index score | 0.4167210932677909 |
| Adjusted Mutual Information Score | 0.5431756896373161 |

| Euclidean metric | |
|---|---|
| Homogeneity score | 0.008428056566878686 |
| Completeness score | 0.008490790660743977 |
| V_measure score | 0.008459307306639803 |
| Adjusted Rand Index score | 0.0014762057832842827 |
| Adjusted Mutual Information Score | 0.00524781170211466 |

# 12  Question 12

Using n_components = 5 as determined above, we get the below contingency matrices – Figure 12 corresponds to the Euclidean metric, and Figure 13 corresponds to the cosine metric.
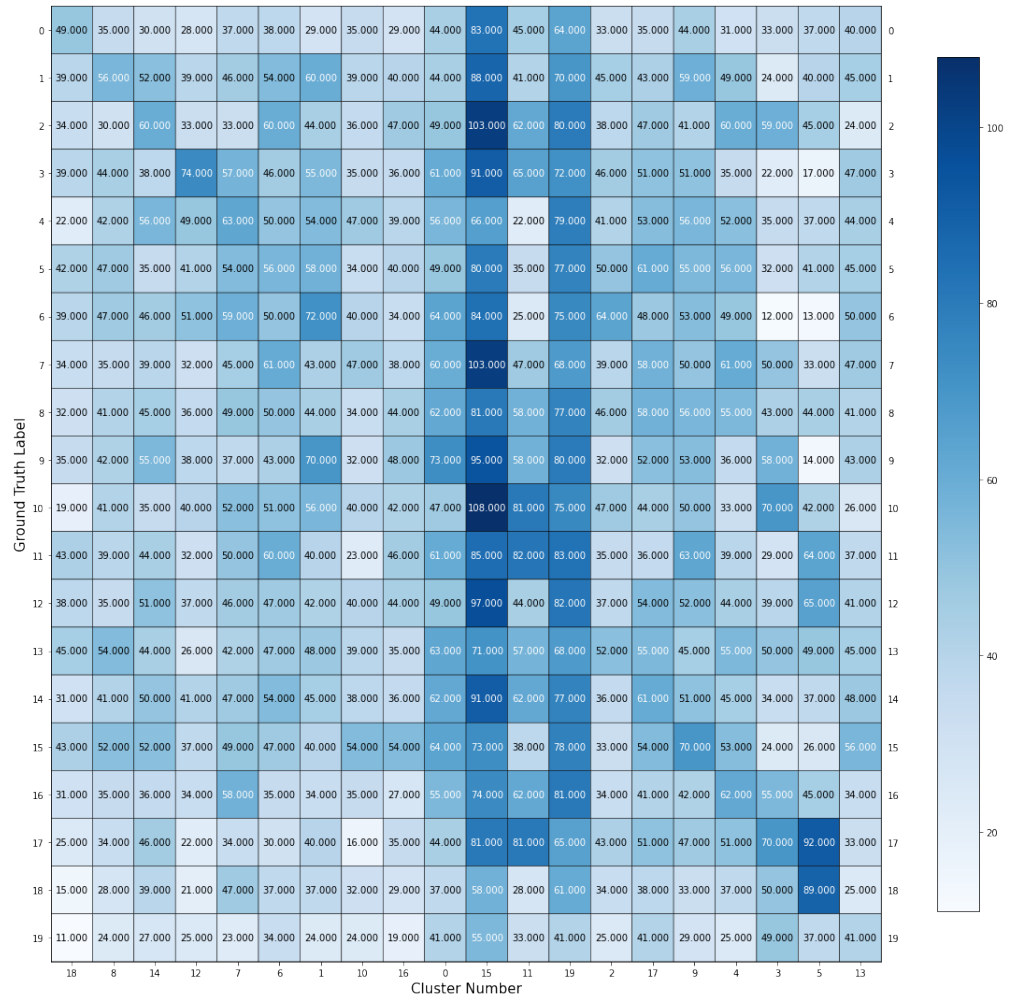
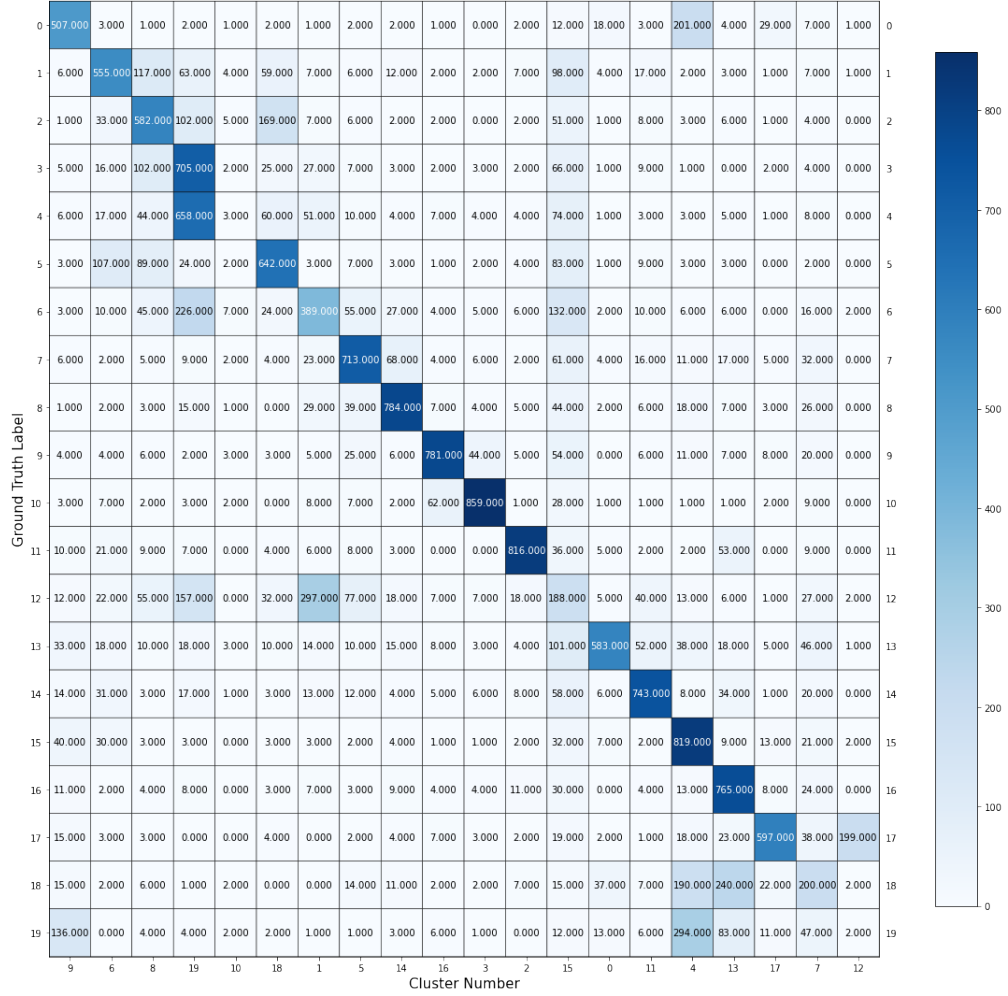Figure 12: Confusion matrix for UMAP-reduced Kmeans with Euclidean metric

Figure 13: Confusion matrix for UMAP-reduced Kmeans with cosine metric

The cosine matrix is obviously much better, which corresponds with the greatly improved scores using this metric.

Examining the cosine metric more closely, we see that most categories were well identified, but some were more likely to be confused with each other. Tracing the connections, we can see that points with ground truth label 19 were likely to be identified with points labeled 0 or 15, and points from label 12 were misidentified as 1 or 6. In the original dataset, this means that points from `talk.religion.misc` were identified as `alt.atheism` or `soc.religion.christian`, and points from `sci.electronics` were identified as `comp.graphics` or `misc.forsale`. The religion groups form an intuitive grouping, and the confusion between the latter three groups perhaps indicates that selling hardware and electronics was

12

common on these newsgroups.

# 13    Question 13

Using the best hyperparameter combination from Q11/12, n_components=5 and the cosine metric, we reduced the Tfidf matrix with UMAP. The scores for agglomerative clustering using 'ward' and 'single' linkage on this data are shown below:

| Ward linkage | |
|---|---|
| Homogeneity score | 0.524473300498607 |
| Completeness score | 0.5479912976714623 |
| V_measure score | 0.53597443687123 |
| Adjusted Rand Index score | 0.3940501791969448 |
| Adjusted Mutual Information Score | 0.5344316766896055 |

| Single linkage | |
|---|---|
| Homogeneity score | 0.022917457355135236 |
| Completeness score | 0.3464156710750703 |
| V_measure score | 0.04299081646294147 |
| Adjusted Rand Index score | 0.00046761208355549543 |
| Adjusted Mutual Information Score | 0.03748616963750802 |

It's clear that ward linkage performs significantly better on the dataset. Single linkage uses the minimum distance between all points in each cluster to determine whether to merge them or not, so it's a lot more sensitive to outliers and tends to result in uneven cluster sizes. Ward linkage on the other hand, minimizes the variance between pairs of clusters, which can produce more regular cluster sizes and hence better results for our particular dataset.

# 14    Question 14

In this section we investigate the DBSCAN and HDBSCAN clustering algorithms. We kept UMAP with n_components=5 and the cosine metric as the dimensionality reduction algorithm. The HDBSCAN min_cluster_size was set to 100, and other hyperparameters were tuned to improve the performance of each of these algorithms. We performed hyperparameter tuning on DBSCAN and HDBSCAN for epsilon values of [0.01,0.05,0.1,0.3,0.5,0.7,0.9,1], metric types of [Euclidean, cosine, Manhattan] (if applicable), and minimum sample values of [10, 20, 50, 70, 100]. The following table depicts the scores for each clustering algorithm for the hyperparameter combination that resulted in the best average score across all 5 clustering metrics:

| DBSCAN: Epsilon = 0.5, Metric = Manhattan, Minimum samples = 10 | |
|---|---|
| Homogeneity score | 0.5407607595568981 |
| Completeness score | 0.45471096739764455 |
| V_measure score | 0.49401673890037173 |
| Adjusted Rand Index score | 0.2967260971944622 |
| Adjusted Mutual Information Score | 0.47920988791262803 |

| HDBSCAN: Epsilon = 0.3, Metric = Euclidean, Minimum samples = 10 | |
|---|---|
| Homogeneity score | 0.43889528158878366 |
| Completeness score | 0.5288819318072416 |
| V_measure score | 0.47970500064413385 |
| Adjusted Rand Index score | 0.2565568616115164 |
| Adjusted Mutual Information Score | 0.4780526456871693 |

# 15   Question 15

DBSCAN was the better performing clustering algorithm. The associated contingency matrix is shown below:
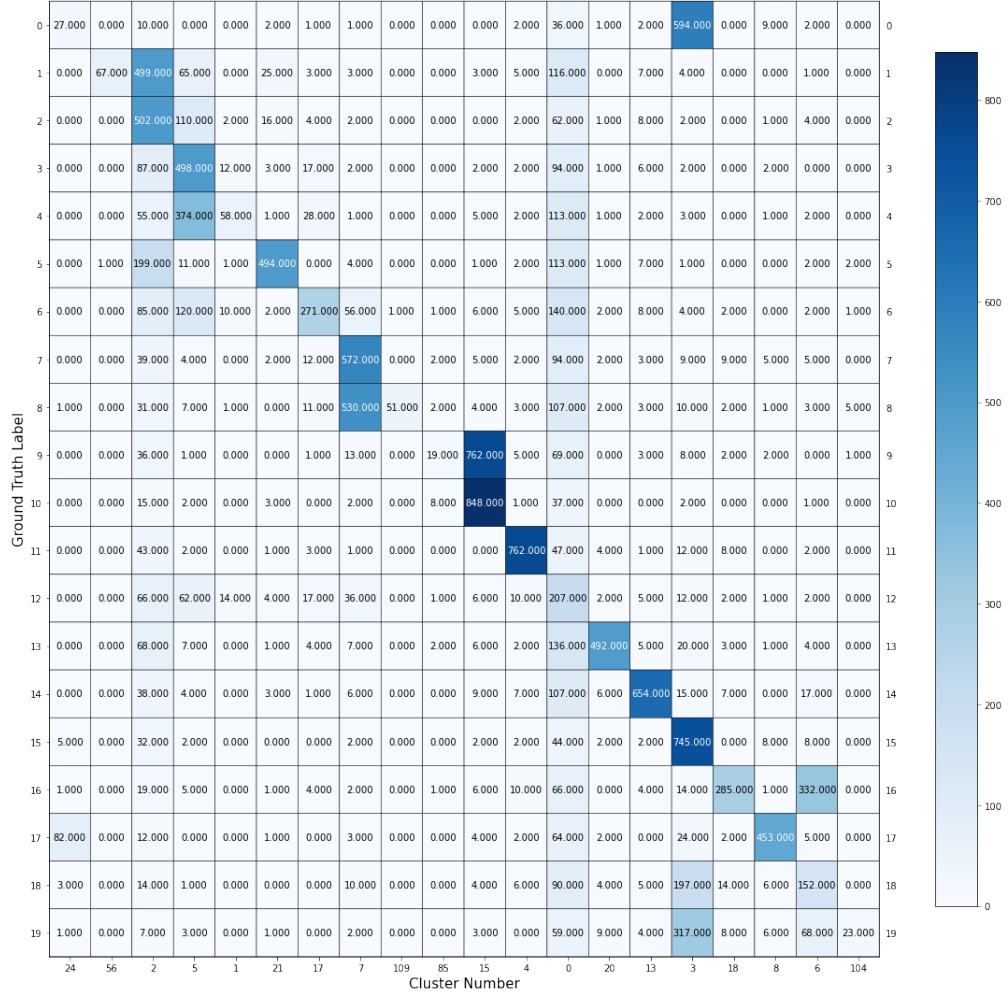
Figure 14: Contingency Matrix for DBSCAN on 20 Class Data

The contingency matrix indicates that there must be at least 96 clusters since that's the largest number shown on the x-axis ticks. By looking at the number of unique labels in the trained DBSCAN model, we are able to determine that there are actually 168 clusters formed by the algorithm (excluding the -1 labels). This number is high because of our low value for the mimimum samples hyperparameter, so most of these cluters contain few data points. However, most of the data points were captured by the larger clusters shown in the contingency matrix, which resulted in acceptable clustering metric scores.

A label of -1 indicates that the data point was classified as noise, and was not assigned to any cluster. This cluster label is represented as cluster number 0 in the contingency matrix. We can tell because this cluster contains many

data points from each of the 20 classes.

# 16 Question 16

We used the "read_csv" function of the Pandas library to acquire the BBC dataset as a dataframe. After converting the written class labels to numerical target values for each document, we then used our established pipeline to vectorize (excluding stopwords/punctuation) then normalize the data and analyze the dimensions of the resulting TF-IDF matrix. As with the 20 Newsgroups dataset, we then reduced dimensionality with UMAP and searched for the optimal parameter configuration for multiple clustering algorithms. Starting with KMeans, we tested how our five aforementioned accuracy metrics vary with the number of principle components used in clustering for both Euclidean (Figure 15) and Cosine distance metrics (Figure 16). The best combination was n_components = 10 with the Euclidean distance metric.
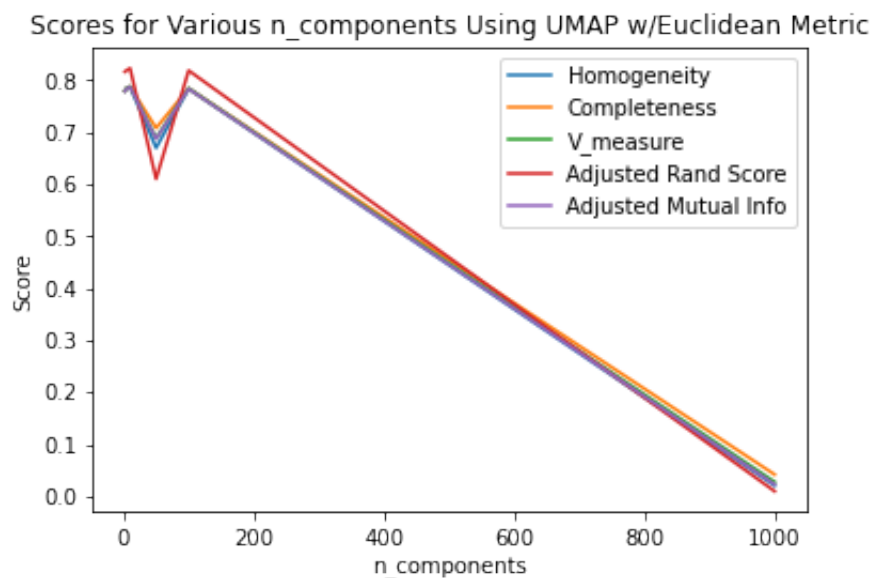


Figure 15: Variation in accuracy metrics by number of components for Kmeans Clustering and Euclidean distance
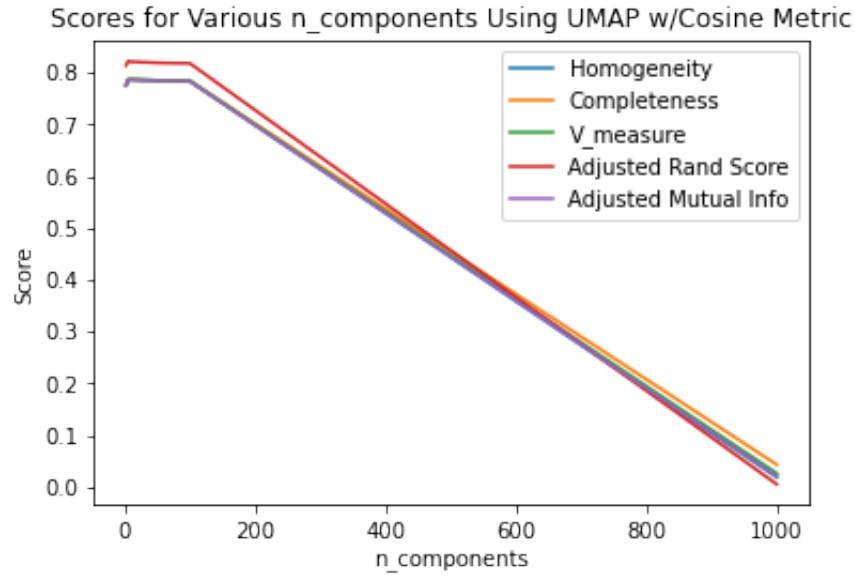
Figure 16: Variation in accuracy metrics by number of components for Kmeans Clustering and Cosine distance

We then tested the accuracy of the DBSCAN clustering algorithm for both distance metrics. We cycled through values from .01 to 1 for epsillon and from 1 to 50 for the minimum number of points per cluster to find the best combination of hyperparameters for our analysis. We used 3-dimensional surface plots to visualize how these parameters in concert affected our five accuracy metrics and the estimated number of clusters generated by DBSCAN (Figure 17).
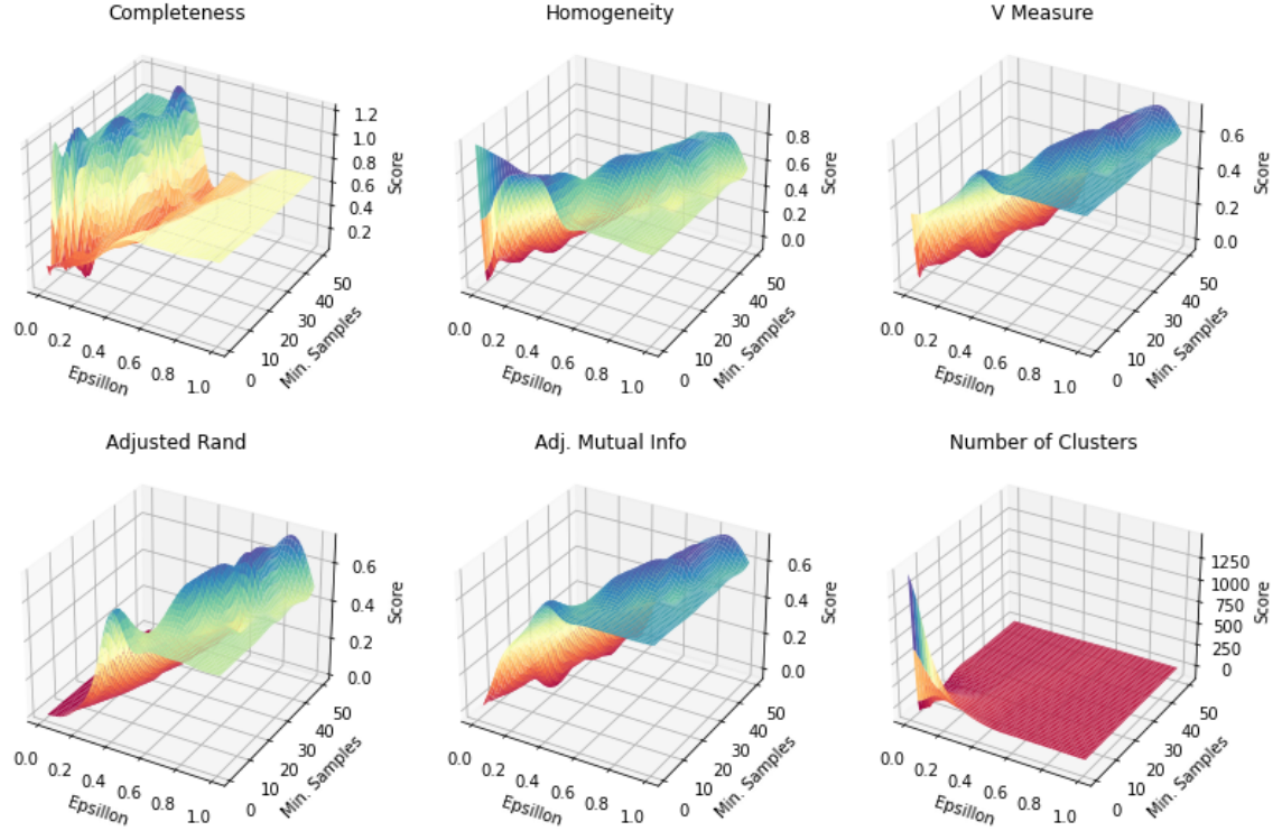
Figure 17: Variation in accuracy metrics corresponding to changes in epsilon and minimum cluster metrics for DBSCAN clustering of UMAP reduced data with Cosine distance metric (Euclidean metric results visually similar)

Avoiding the low values of epsilon and minimum samples that resulted in excessively high numbers of estimated clusters, we found the combination of parameters that maximized the average of the accuracy metrics to be 0.8 for epsilon and 35 for minimum samples. We then compared the results of each of our protocols (Kmeans with Euclidean/Cosine, DBSCAN with Euclidean/Cosine) to find the best reduction and clustering technique for our data. We report the five accuracy metrics for each protocol below:

| Accuracy Metrics for DBSCAN, Metric = Euclidean | |
|---|---|
| Homogeneity score | 0.7896799513244787 |
| Completeness score | 0.667782538810624 |
| V_measure score | 0.7236336939202158 |
| Adjusted Rand Index score | 0.7108732847743878 |
| Adjusted Mutual Information Score | 0.7221232251101021 |

| Accuracy Metrics for DBSCAN, Metric = Cosine | |
|---|---|
| Homogeneity score | 0.7896799513244787 |
| Completeness score | 0.667782538810624 |
| V_measure score | 0.7236336939202158 |
| Adjusted Rand Index score | 0.7108732847743878 |
| Adjusted Mutual Information Score | 0.7221232251101021 |

| Accuracy Metrics for KMeans, Metric = Euclidean | |
|---|---|
| Homogeneity score | 0.7874360944663693 |
| Completeness score | 0.7874929063651126 |
| V_measure score | 0.7874644993910623 |
| Adjusted Rand Index score | 0.8229763601775667 |
| Adjusted Mutual Information Score | 0.786746935945491 |

| Accuracy Metrics for KMeans, Metric = Cosine | |
|---|---|
| Homogeneity score | 0.7863911146778318 |
| Completeness score | 0.7863053396249526 |
| V_measure score | 0.7863482248123015 |
| Adjusted Rand Index score | 0.8206748429226303 |
| Adjusted Mutual Information Score | 0.785626959948192 |

The best scoring protocol was Kmeans with Euclidean distance metric and n_components set to 10. The contingency matrix for this clustering is shown in Figure 18 below
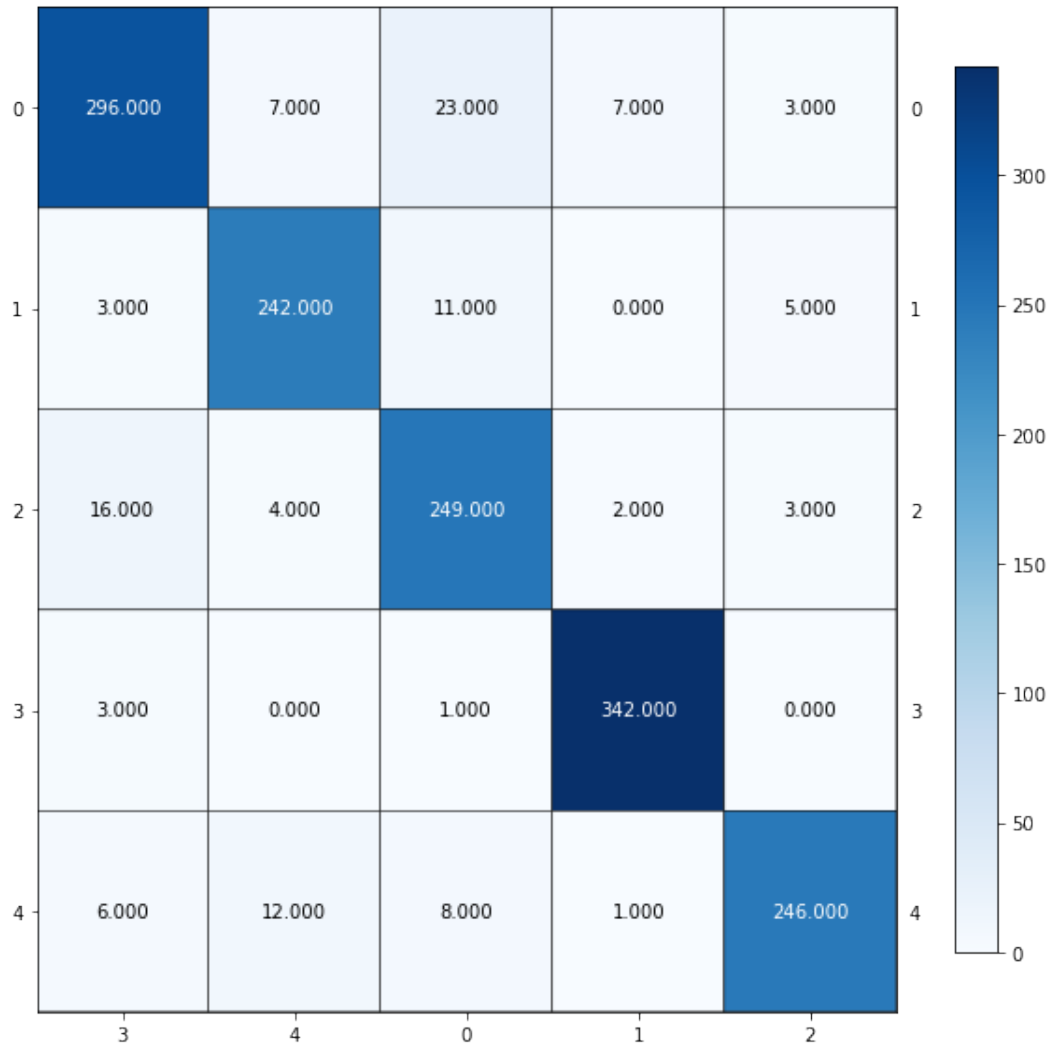
Figure 18: Contingency Matrix for optimal clustering of BBC News data

Our pipeline and parameter selections produce a fairly good clustering through both DBSCAN and KMeans. One factor that may affect the accuracy of our analysis is the relatively small size of the BBC dataset when compared to the 20 Newsgroups dataset. With more data, our clustering would likely become more accurate.