# ECE219 - Project 2 Report

# Clustering

2/12/18

Jonathan Hurwitz (804258351)

Peter Kim (204271299)

# 1.

Here we transform our newsgroup documents into TF-IDF vector representations. We use min_df=3 and filter out stop words. Stemming was not involved. The final dimension of our matrix was 7882 documents X 27768 tokens.

# 2.

a) The contingency matrix after applying k-means with r=2 on the previous vectorized dataset was as follows:

| 4 | 3889 |
|---|---|
| 1728 | 2251 |

b) The clustering metrics were as follows:
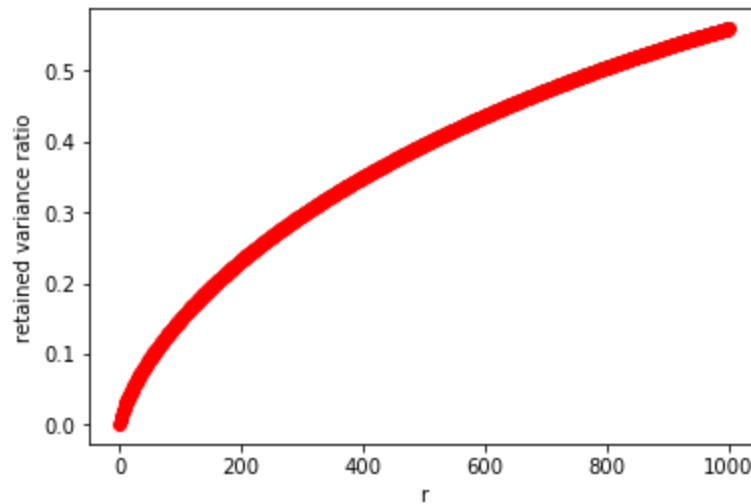
| Homogeneity | 0.253 |
|---|---|
| Completeness | 0.335 |
| V-measure | 0.288 |
| Adjusted Rand-Index | 0.181 |
| Adjusted Mutual Info | 0.253 |

The initial TF-IDF representations are high dimensional and sparse. The clustering performance given high dimensional and sparse data is poor. This is expected as Euclidean distances, which the k-means algorithm attempts to minimize every iteration, are essentially reduced to be the same for almost all pairs of examples in high-dimensional settings. This, therefore, renders k-means as an ineffective clustering technique for high-dimensions due to the curse of dimensionality.

# 3.

In this section, we preprocess the TF-IDF matrix by reducing the number of feature dimensions by applying both Latent Semantic Indexing (LSI) and Non-negative Matrix Factorization (NMF).
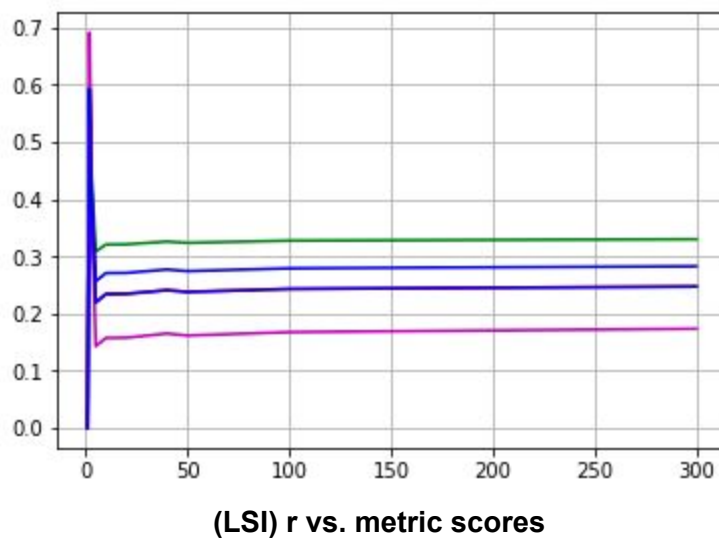
ai) We first attempted to find the most effective dimension to reduce to through inspecting change in the retained variance ratio when removing the top singular values. For the top 1000 principle components, we show the retained variance ratio in the following graph:
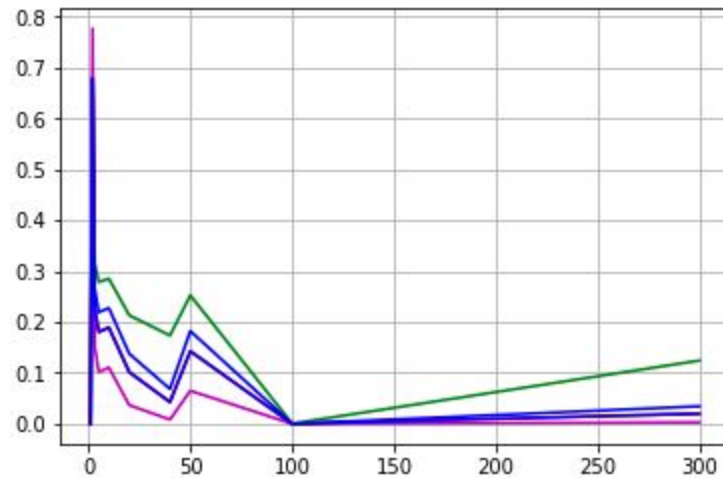


At r=1000 principal components, we retain around 60% of the total variance.

aii) Afterwards, we used the two methods of LSI and NMF and try r values in the range [1,2,3,5,10,20,50,100,300] and determined the clustering performance with respect to the five metrics mentioned above.

When reducing with LSI and NMF, we get the following performance scores across our range of r.



**(LSI) r vs. metric scores**

**(NMF) r vs. metric scores**

The non-monotonicity of the two above graphs can be explained as such: as we increase the number of dimensions r, more information is encapsulated in the feature representation just by nature of having a higher percentage of the variance due to more principal components. However, K-Means suffers from the curse of dimensionality, so as the dimensions increase, Euclidean distance becomes a worse and worse distance metric and will perform poorly. Thus, we are trading off between more information and an inherent shortcoming in the K-Means method. This is why it's possible for kmeans to perform better with lower dimension (such as r=100) than with higher dimension (such as r=300 or r=1000).

From the above graphs as well as the metrics shown in the table below, it is clear that **the best r-value for both LSI and NMF is r=2**. As mentioned before, we are trading off between incorporating more information and the high dimensionality shortcomings of K-Means. It makes sense that a low dimension value such as 2 would be effective for K-Means, due to the aforementioned curse of dimensionality.

The contingency matrices and corresponding metrics for all cases are summarized in the following table:

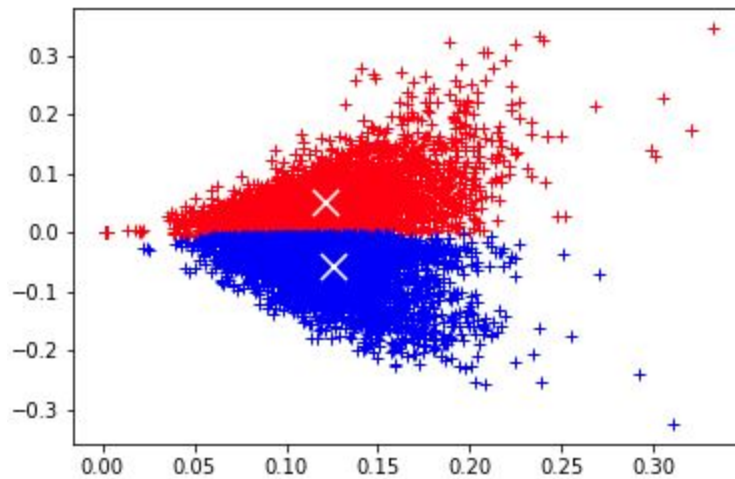| r-value | LSI | NMF |
|---|---|---|
| 1 | homogeneity=0.00030<br>completeness=0.00030<br>v-measure=0.00030<br>adj rand index=0.00034<br>adj mutual info=0.00021<br>[[1703 2200]<br>[1656 2323]] | homogeneity=0.00030<br>completeness=0.00030<br>v-measure=0.00030<br>adj rand index=0.00034<br>adj mutual info=0.00021<br><br>[[2200 1703] |

| | | |
|---|---|---|
| | | [2323 1656]] |
| 2<br>Best value for LSI<br>Best value for NMF | **homogeneity=0.59263**<br>**completeness=0.59424**<br>**v-measure=0.59343**<br>**adj rand index=0.69180**<br>**adj mutual info=0.59259**<br><br>**[[3713  190]**<br>**[ 473 3506]]** | **homogeneity=0.67905**<br>**completeness=0.68013**<br>**v-measure=0.67959**<br>**adj rand index=0.77702**<br>**adj mutual info=0.67902**<br><br>[[3594  309]<br>[ 158 3821]] |
| 3 | homogeneity=0.40097<br>completeness=0.43866<br>v-measure=0.41897<br>adj rand index=0.39656<br>adj mutual info=0.40091<br><br>[[3866   37]<br>[1422 2557]] | homogeneity=0.22934<br>completeness=0.31648<br>v-measure=0.26596<br>adj rand index=0.15280<br>adj mutual info=0.22927<br><br>[[3899    4]<br>[2396 1583]] |
| 5 | homogeneity=0.21960<br>completeness=0.30838<br>v-measure=0.25653<br>adj rand index=0.14284<br>adj mutual info=0.21953<br><br>[[3898    5]<br>[2446 1533]] | homogeneity=0.18063<br>completeness=0.27871<br>v-measure=0.21920<br>adj rand index=0.10196<br>adj mutual info=0.18056<br><br>[[3898    5]<br>[2677 1302]] |
| 10 | homogeneity=0.23427<br>completeness=0.32095<br>v-measure=0.27084<br>adj rand index=0.15740<br>adj mutual info=0.23420<br><br>[[3900    3]<br>[2374 1605]] | homogeneity=0.18920<br>completeness=0.28527<br>v-measure=0.22751<br>adj rand index=0.11056<br>adj mutual info=0.18913<br><br>[[3898    5]<br>[2625 1354]] |
| 20 | homogeneity=0.23463<br>completeness=0.32122 | homogeneity=0.10163<br>completeness=0.21326<br>v-measure=0.13766 |

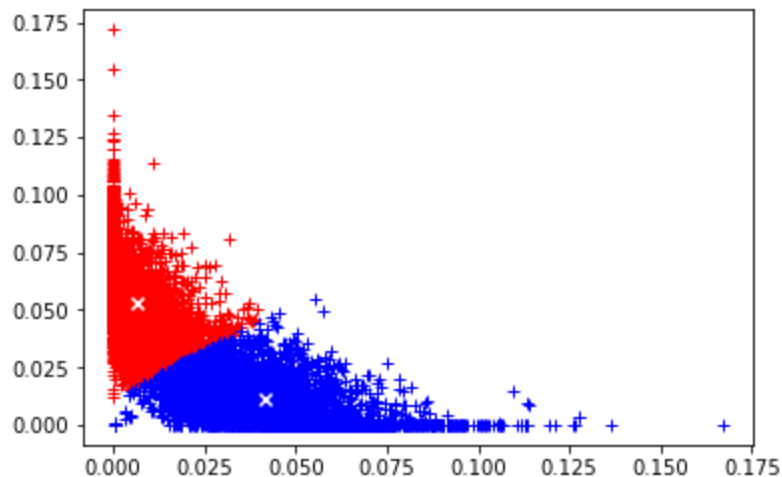| | | |
|---|---|---|
| | v-measure=0.27118<br>adj rand index=0.15780<br>adj mutual info=0.23455<br><br>[[   3 3900]<br> [1607 2372]] | adj rand index=0.03660<br>adj mutual info=0.10155<br><br>[[   7 3896]<br> [ 800 3179]] |
| 50 | homogeneity=0.23801<br>completeness=0.32377<br>v-measure=0.27434<br>adj rand index=0.16165<br>adj mutual info=0.23794<br><br>[[3900    3]<br> [2353 1626]] | homogeneity=0.14268<br>completeness=0.25245<br>v-measure=0.18232<br>adj rand index=0.06493<br>adj mutual info=0.14260<br><br>[[   2 3901]<br> [1045 2934]] |
| 100 | homogeneity=0.24321<br>completeness=0.32768<br>v-measure=0.27920<br>adj rand index=0.16763<br>adj mutual info=0.24314<br><br>[[3900    3]<br> [2324 1655]] | homogeneity=0.00002<br>completeness=0.00045<br>v-measure=0.00003<br>adj rand index=0.00001<br>adj mutual info=-0.00008<br><br>[[3887   16]<br> [3965   14]] |
| 300 | homogeneity=0.24740<br>completeness=0.33013<br>v-measure=0.28284<br>adj rand index=0.17350<br>adj mutual info=0.24733<br><br>[[   4 3899]<br> [1684 2295]] | homogeneity=0.01995<br>completeness=0.12427<br>v-measure=0.03438<br>adj rand index=0.00283<br>adj mutual info=0.01986<br><br>[[3723  180]<br> [3974    5]] |

**4.**

**a)** The best clustering results for LSI and NMF were obtained for r=2. This value was used in the dimensionality reduction. First, we visualized the performance of the best clustering results by projecting the final data vectors onto a 2 dimensional place and color coding the two classes.



LSI-reduced clustering result with r=2
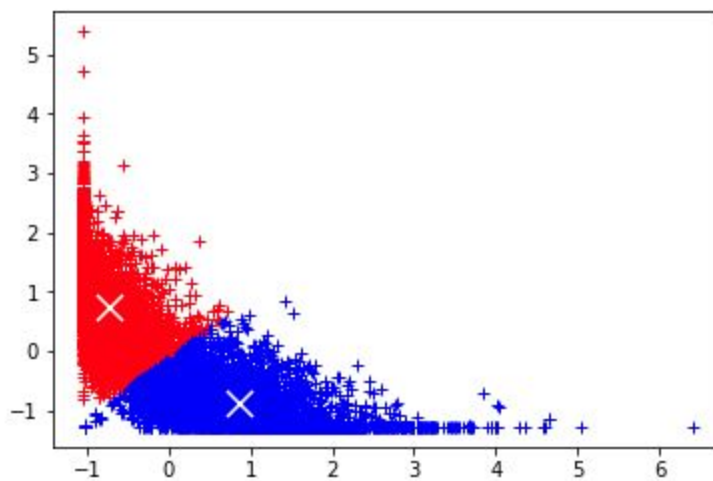


NMF-reduced clustering result with r=2

We are content with these clustering results as the two different types of classes are generally centered around their respective centroids (marked with a white X) in a cohesive manner.

**b)** The next step was to try several different methods to see if they increase clustering performance for r=2 NMF reduced data. The methods and their results are described below:

   1)  Normalizing the features such that each feature has unit variance.

      The table below shows a comparison between the baseline K-Means + NMF (r=2) results and the normalized NMF (r=2) results. NMF with normalization results in marginally better homogeneity, completeness, v-measure, and adjusted mutual information. However, adjusted random index is marginally lower.

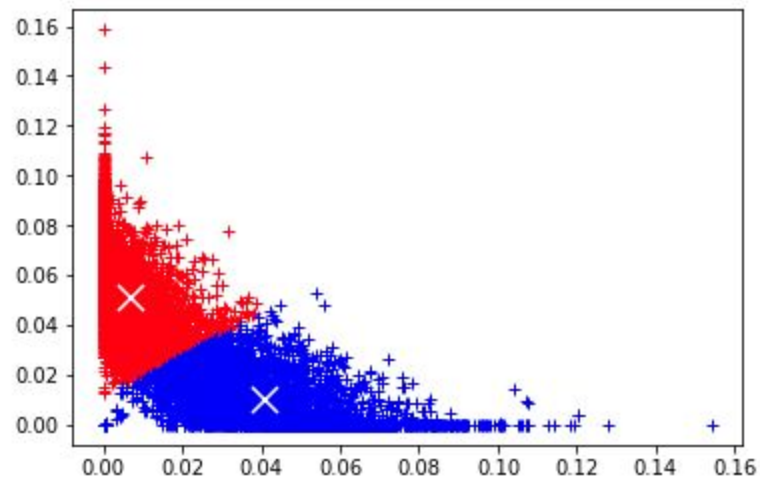| NMF (r=2) baseline performance | NMF (r=2) with normalization performance |
|---|---|
| **homogeneity=0.67905**<br> **completeness=0.68013**<br> **v-measure=0.67959**<br> **adj rand index=0.77702**<br> **adj mutual info=0.67902**<br><br>[[3594  309]<br> [ 158 3821]] | **homogeneity=0.68280**<br>**completeness=0.68564**<br>**v-measure=0.68422**<br> **adj rand index=0.77344**<br> **adj mutual info=0.68277**<br><br>[[ 369, 3534],<br>     [3873,  106]] |



**NMF (r=2) with Normalization**

2) Applying non-linear transformation (log transformation).

The table below shows a comparison between the baseline performance and K-Means after a log transformation is applied. All of the metrics decrease.

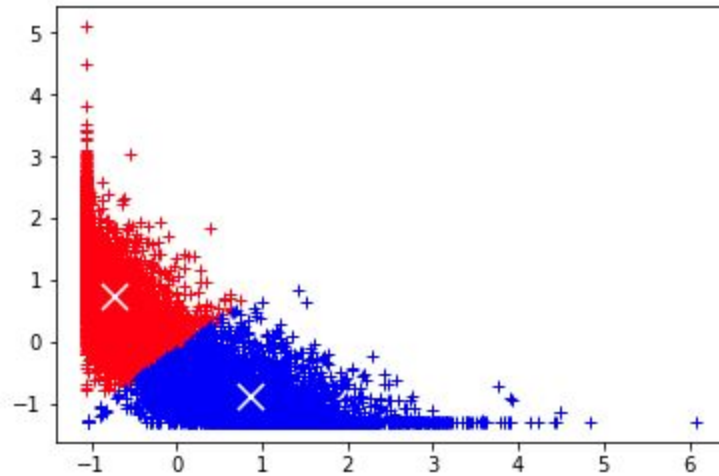| NMF (r=2) baseline performance | NMF (r=2) with log transformation performance |
|---|---|
| **homogeneity=0.67905**<br> **completeness=0.68013**<br> **v-measure=0.67959**<br> **adj rand index=0.77702**<br> **adj mutual info=0.67902**<br><br>[[3594  309]<br> [ 158 3821]] | **homogeneity=0.67505**<br>**completeness=0.67641**<br>**v-measure=0.67573**<br> **adj rand index=0.77255**<br> **adj mutual info=0.67503**<br><br>[[ 325, 3578],<br> [3827,  152]] |

|  |  |
|---|---|
|  |  |



**NMF (r=2) with log transformation**

3) Applying log transformation then normalization.

NMF with log transformation then normalization showed an improvement across all metrics except for adjusted rand index.

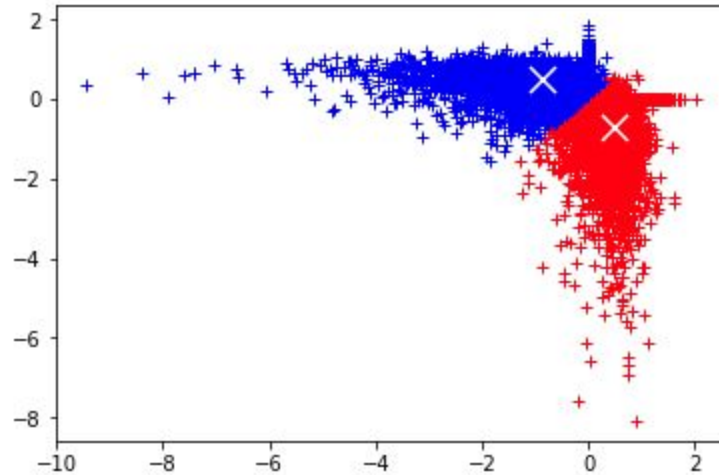| NMF (r=2) baseline performance | NMF (r=2) with log transformation then normalization performance |
|---|---|
| **homogeneity=0.67905**<br> **completeness=0.68013**<br> **v-measure=0.67959**<br> **adj rand index=0.77702**<br> **adj mutual info=0.67902**<br><br>[[3594  309]<br> [ 158 3821]] | **homogeneity=0.68369**<br>**completeness=0.68649**<br>**v-measure=0.68509**<br> **adj rand index=0.77434**<br> **adj mutual info=0.68366**<br><br>[[ 367, 3536],<br>  [3873,  106]] |

**NMF (r=2) with log transformation then normalization**

4) Applying normalization then log transformation.

NMF with normalization and then log transformation showed a decrease in performance across all categories.

| NMF (r=2) baseline performance | NMF (r=2) with normalization then log transformation performance |
|---|---|
| **homogeneity=0.67905**<br> **completeness=0.68013**<br> **v-measure=0.67959**<br> **adj rand index=0.77702**<br> **adj mutual info=0.67902**<br><br>[[3594  309]<br> [ 158 3821]] | **homogeneity=0.65169**<br>**completeness=0.65196**<br>**v-measure=0.65183**<br> **adj rand index=0.75525**<br> **adj mutual info=0.65166**<br><br>[[3613,  290],<br>   [ 226, 3753]] |

**NMF (r=2) with normalization then log transformation**

A non-linear transformation such as a logarithmic transformation can increase the relationship between independent features, which are the case for our TFIDF matrix. This can add linearity to our data and thus, could increase the clustering results. However, in our case of applying logarithmic transformation to the NMF-reduced dataset, we empirically show that the performance was actually hurt by such a nonlinear transformation (~1% decrease).

# 5.

This section involved clustering of 20 categories rather than just 2. All documents were included in the data matrix. The same parameters as in part 1 were used with regards to stemming and min_df. We ran K-Means with varying values of r for both LSI and NMF dimensionality reduction to find an optimal value for each. The results from this sweep are summarized in the table below.

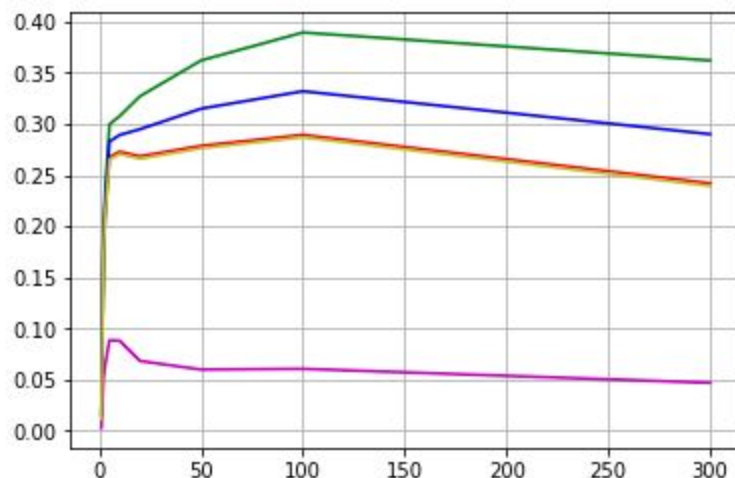| r-value | LSI | NMF |
|---|---|---|
| 1 | homogeneity=0.01516<br> completeness=0.01648<br> v-measure=0.01579<br> adj rand index=0.00313<br> adj mutual info=0.01196 | homogeneity=0.01546<br> completeness=0.01666<br> v-measure=0.01604<br> adj rand index=0.00316<br> adj mutual info=0.01226 |
| 2 | homogeneity=0.17440<br> completeness=0.18632 | homogeneity=0.16135<br> completeness=0.17094<br> v-measure=0.16601 |

| | | |
|---|---|---|
| | v-measure=0.18016<br>adj rand index=0.05036<br>adj mutual info=0.17173 | adj rand index=0.04657<br>adj mutual info=0.15863 |
| 3 | homogeneity=0.21548<br>completeness=0.23676<br>v-measure=0.22562<br>adj rand index=0.06647<br>adj mutual info=0.21293 | homogeneity=0.18779<br>completeness=0.21056<br>v-measure=0.19852<br>adj rand index=0.05472<br>adj mutual info=0.18515 |
| 5 | homogeneity=0.26561<br>completeness=0.29715<br>v-measure=0.28050<br>adj rand index=0.08719<br>adj mutual info=0.26323 | homogeneity=0.21922<br>completeness=0.26075<br>v-measure=0.23819<br>adj rand index=0.06541<br>adj mutual info=0.21666 |
| 10<br>Best value for NMF | homogeneity=0.27704<br>completeness=0.31739<br>v-measure=0.29584<br>adj rand index=0.08653<br>adj mutual info=0.27469 | **homogeneity=0.26157**<br>**completeness=0.31635**<br>**v-measure=0.28636**<br>**adj rand index=0.07140**<br>**adj mutual info=0.25917** |
| 20 | homogeneity=0.28841<br>completeness=0.35187<br>v-measure=0.31699<br>adj rand index=0.07260<br>adj mutual info=0.28609 | homogeneity=0.22945<br>completeness=0.29031<br>v-measure=0.25632<br>adj rand index=0.04733<br>adj mutual info=0.22694 |
| 50 | homogeneity=0.29119<br>completeness=0.38053<br>v-measure=0.32992<br>adj rand index=0.06287<br>adj mutual info=0.28888 | homogeneity=0.16791<br>completeness=0.24519<br>v-measure=0.19932<br>adj rand index=0.02413<br>adj mutual info=0.16520 |
| 100<br>Best value for LSI | **homogeneity=0.29706**<br>**completeness=0.38877**<br>**v-measure=0.33679** | homogeneity=0.06814<br>completeness=0.11786<br>v-measure=0.08636 |

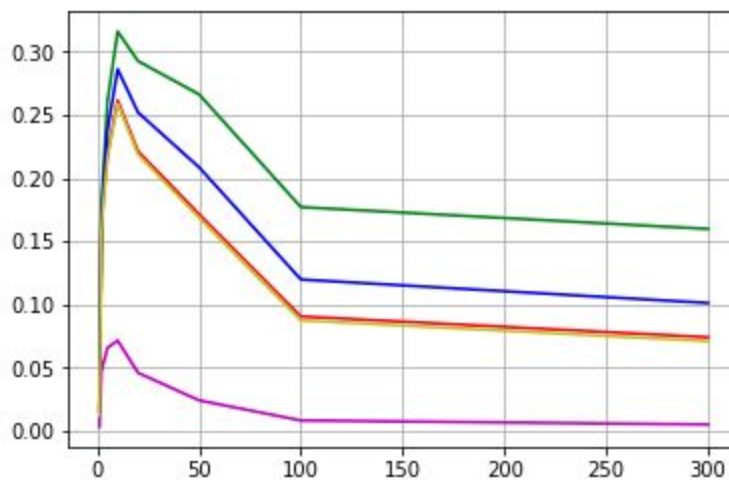| | adj rand index=0.06609<br>adj mutual info=0.29477 | adj rand index=0.00549<br>adj mutual info=0.06504 |
| --- | --- | --- |
| 300 | homogeneity=0.28491<br>completeness=0.39868<br>v-measure=0.33233<br>adj rand index=0.06989<br>adj mutual info=0.28256 | homogeneity=0.07278<br>completeness=0.12840<br>v-measure=0.09290<br>adj rand index=0.00654<br>adj mutual info=0.06973 |

The best r-value for LSI was r=100. The best value for NMF was r=10.

The following plots show the performance for both LSI and NMF as r-value increases. The y-axis is units of whatever metric is being displayed on a scale from 0 to 1. The x-axis represents r-value. The legend is as follows:

| Color | Metric |
| --- | --- |
| red | Homogeneity |
| green | Completeness score |
| blue | V-measure |
| magenta | Adjusted rand-index |
| yellow | Adjusted mutual information |



**(LSI) metric scores vs. r-value**

**(NMF) metric scores vs. r-value**

The optimal r-values for LSI and NMF (100 and 10 respectively) were used in the transformation tests performed below. The same four methods were tested:

1) Normalizing the features such that each feature has unit variance.

For LSI, applying normalization worsened all of the metrics. Normalization improved adjusted rand index for NMF, but worsened all of the other metrics.

| LSI baseline (r=100) | LSI normalized (r=100) | NMF baseline (r=10) | NMF normalized (r=10) |
|---|---|---|---|
| homogeneity=0.29706 completeness=0.38877 v-measure=0.33679 adj rand index=0.06609 adj mutual info=0.29477 | homogeneity=0.19519 completeness=0.27588 v-measure=0.22863 adj rand index=0.02468 adj mutual info=0.19257 | homogeneity=0.26157 completeness=0.31635 v-measure=0.28636 adj rand index=0.07140 adj mutual info=0.25917 | homogeneity=0.25016 completeness=0.29413 v-measure=0.27036 adj rand index=0.08328 adj mutual info=0.24772 |

2) Applying non-linear transformation (log transformation).

Applying the log transform to LSI lowered all metric values. For NMF, homogeneity, v-measure, adjusted rand index, and adjusted mutual information improved while completeness marginally decreased.

| LSI baseline (r=100) | LSI log (r=100) | NMF baseline (r=10) | NMF log (r=10) |
|---|---|---|---|
| homogeneity=0.29706 completeness=0.38877 v-measure=0.33679 | homogeneity=0.28039 completeness=0.37745 v-measure=0.32176 | homogeneity=0.26157 completeness=0.31635 v-measure=0.28636 adj rand index=0.07140 | homogeneity=0.26394 completeness=0.31519 v-measure=0.28729 adj rand index=0.07274 |

| | | | |
|---|---|---|---|
| adj rand index=0.06609<br>adj mutual info=0.29477 | adj rand index=0.06485<br>adj mutual info=0.27803 | adj mutual info=0.25917 | adj mutual info=0.26154 |

3) Applying log transformation then normalization.

Applying log transformation and then normalization worsened all of the metrics across the board for both LSI and NMF, with the exception of adjusted rand index for NMF.

| LSI baseline (r=100) | LSI log + scale (r=100) | NMF baseline (r=10) | NMF log + scale (r=10) |
|---|---|---|---|
| homogeneity=0.29706<br>completeness=0.38877<br>v-measure=0.33679<br>adj rand index=0.06609<br>adj mutual info=0.29477 | homogeneity=0.20665<br>completeness=0.30849<br>v-measure=0.24750<br>adj rand index=0.03525<br>adj mutual info=0.20404 | homogeneity=0.26157<br>completeness=0.31635<br>v-measure=0.28636<br>adj rand index=0.07140<br>adj mutual info=0.25917 | homogeneity=0.24761<br>completeness=0.28635<br>v-measure=0.26558<br>adj rand index=0.07678<br>adj mutual info=0.24516 |

4) Applying normalization then log transformation.

Applying feature scaling before the log transformation resulted in score improvements for all metrics across the board. **Empirically, this was the best method, so we can conclude that feature scaling should be applied before any type of non-linear transformation.**

| LSI baseline (r=100) | LSI scale + log (r=100) | NMF baseline (r=10) | NMF scale + log (r=10) |
|---|---|---|---|
| homogeneity=0.29706<br>completeness=0.38877<br>v-measure=0.33679<br>adj rand index=0.06609<br>adj mutual info=0.29477 | homogeneity=0.30133<br>completeness=0.38911<br>v-measure=0.34348<br>adj rand index=0.07849<br>adj mutual info=0.30230 | homogeneity=0.26157<br>completeness=0.31635<br>v-measure=0.28636<br>adj rand index=0.07140<br>adj mutual info=0.25917 | homogeneity=0.26450<br>completeness=0.27809<br>v-measure=0.27113<br>adj rand index=0.12484<br>adj mutual info=0.26212 |