# Clustering Wikipedia Articles

**Lane Aasen**
Department of Computer Science
University of Washington
Seattle, WA 98105
`aaasen@cs.washington.edu`

## Abstract

Clustering Wikipedia articles using unsupervised learning techniques including K-Means and Latent Dirichlet Allocation (LDA). Exploration into K-Means including choice of K, K-Means++, the effect of random seeds, and TF-IDF thresholds.

## 1 Dataset

The provided dataset contains 15,903 Wikipedia articles in term frequency inverse document frequency (tf-idf) format. Most of the articles are in English, but some are not. Sorting the documents by language would be an interesting project in itself! There are 10,574 unique words in this dataset. Each document is represented as a sparse vector with one dimension for each word. Stop words have been removed from the dataset, but uncommon words remain.

## 2 K-Means Clustering

My first goal was to implement a basic K-Means clustering algorithm from scratch. The code is available at `https://github.com/aaasen/wiki-cluster`.

### 2.1 Choosing K

#### 2.1.1 Minimizing Distortion

Given $K$ clusters $C_1, C_2, ..., C_K$ where each cluster is a set of document vectors and $\mu_i$ is the centroid of $C_i$, the total distortion is defined as follows:

$$\sum_{i=1}^{K} \sum_{d \in C_i} ||d - \mu_i||^2$$

To minimize the distortion, we could set $K$ equal to the number of documents, but then the clusters would be meaningless. We want to choose a $K$ with low distortion that also results in interpretable clusters. Figure 1 shows a plot of $K$ versus total distortion. When $1 \leq K \leq 16$, adding additional clusters has a large impact on the distortion, but once $K > 16$, adding additional clusters has little impact on the distortion. From this alone, it makes sense to set $K = 16$ since it provides a good balance of distortion and interpretability.
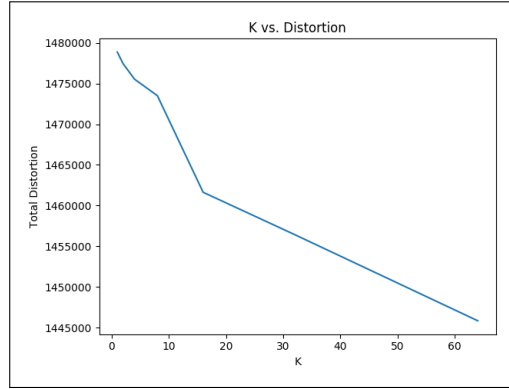
Figure 1: $K$ versus total distortion for $K \in \{1, 2, 4, ..., 256\}$

## 2.2 K and Cluster Size

As $K$ increases, the clusters become more sparse. Once $K = 256$, over half of the clusters have only one document, and are essentially useless. When $K = 16$, the median cluster size is 8.5, and the cluster sizes are as follows:

$$[10061, 3013, 1128, 909, 707, 30, 23, 13, 4, 4, 3, 2, 2, 2, 1, 1]$$

Over half of the clusters are very small, and one of the clusters is too large to be interpretable. This indicates that the data has significant outliers and may lack a structure conducive to clustering.
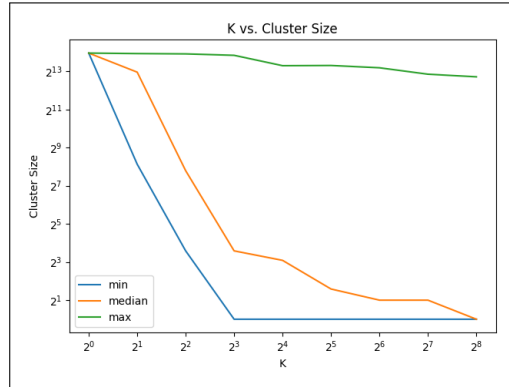


Figure 2: $K$ versus minimum, median, and maximum cluster size for $K \in \{1, 2, 4, ..., 256\}$ with a $log_2$ scale on both axes.

## 2.3 Exploring Clusters

Table 1 shows the clusters with at least 10 documents for K-Means clustering with $K = 16$. The words in each cluster are the dimensions of the centroid with the largest magnitude. The documents shown are those that are closest to the centroid of the cluster.

Overall, the generated clusters make sense, but there are some points of confusion:

- The words that make up cluster 0 have little relation to each other. This cluster contains the majority of the documents.
- Cluster 1 contains churches as well as colleges.

2

- Cluster 3 contains documents related to TV shows and sports because both contain the word "season." Another explanation for this is that sports players appear on Dancing with the Stars.

Table 1: K-Means clusters with $K = 16$ and at least 10 documents.

| Cluster | Size | Words | Documents |
|---|---|---|---|
| 0 | 10061 | females<br>station<br>family<br>located<br>north | mcgillpainquestionnaire<br>historyofthefamily<br>thetussaudsgroup<br>nadiraactress<br>mansfieldsummithighschool |
| 1 | 3013 | church<br>college<br>students<br>published<br>institute | edmondscommunitycollege<br>helderbergcollege<br>oberlincongregationalchurch<br>lundbyoldchurch<br>dioceseoflimerickandkillaloe |
| 2 | 1128 | party<br>served<br>general<br>member<br>senate | partyidentification<br>labourfarmerparty<br>democraticalliancesouthafrica<br>liberaldemocratsitaly<br>christiancreditparty |
| 3 | 909 | season<br>club<br>playing<br>seasons<br>player | dancingwiththestars<br>davidmccracken<br>gilbertcurgenven<br>bjsamsamericanfootball<br>livingstonewalker |
| 4 | 707 | album<br>released<br>songs<br>records<br>rock | thegreatestdaytakethatalbum<br>conflictingemotions<br>primalscream<br>leftbacklp<br>elisamartin |
| 5 | 30 | nba<br>basketball<br>points<br>season<br>seasons | kcjones<br>hakeemolajuwon<br>albertkingbasketball<br>ballstatecardinalsmensbasketball<br>201011southfloridabullsmensbasketballteam |
| 6 | 23 | riots<br>police<br>murder<br>captured<br>robbery | sowetouprising<br>1992losangelesriots<br>nikolaybogolepov<br>josephlamothe<br>jenmi |
| 7 | 13 | congo<br>subtropical<br>republic<br>zambia<br>zimbabwe | republicofcabinda<br>brownrumpedbunting<br>copperbeltprovince<br>leptopelisviridis<br>yellowthroatedpetronia |

## 2.4 K-Means++

K-Means++ is a method for selecting the initial cluster centers for K-Means. In the normal version of K-Means, initial cluster centroids are selected by sampling random documents from the dataset. This can produce suboptimal clusters, increase time to convergence, and increase variability in clustering performance.

K-Means++ solves this by selecting initial centroids one at a time in order to minimize distortion. Each document is given a weight proportional to the distance to the nearest centroid and the next centroid is selected at random using these weights until $K$ centroids have been chosen. This results in better centroids than random sampling and lessens the time until convergence. However, it can

be quite expensive when $K$ is large since the distance from each document to each cluster centroid must be computed at every iteration.

I implemented K-Means++ and expected it to increase clustering performance considerably. However, I found that it actually had a slightly negative impact on distortion and noticeably increased cluster sparsity.

I think that this is because outliers were chosen as initial centroids and the clusters never expanded to include other points.
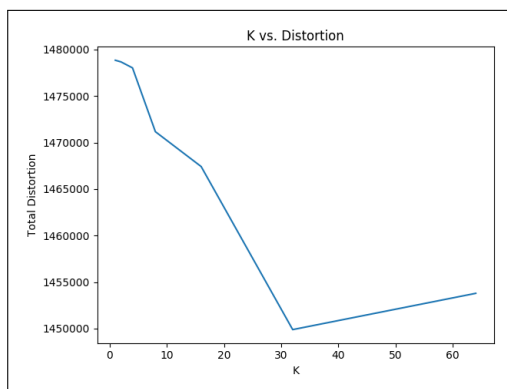


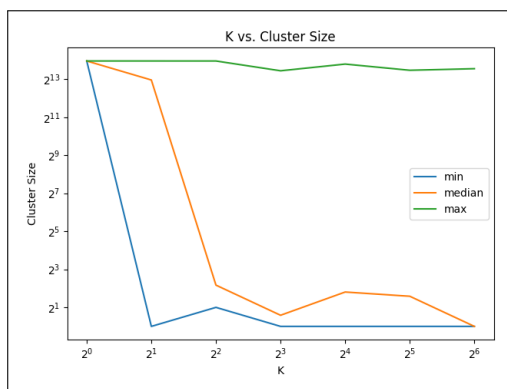Figure 3: $K$ versus total distortion for $K \in \{1, 2, 4, ..., 64\}$ using K-Means++.



Figure 4: $K$ versus minimum, median, and maximum cluster size for $K \in \{1, 2, 4, ..., 64\}$ with a $log_2$ scale on both axes. Using K-Means++.

## 2.5 Fighting Cluster Sparsity

After implementing K-Means++ with no success, it became apparent that there was a problem with the underlying data, and not with the clustering algorithm.

I began to investigate the provided dataset, which was in tf-idf format with no documentation. I found that stop words had been removed from the dataset, but rare words had been left in. Half of the words in the dictionary were used in 16 or fewer datasets (less than 0.1%). Tables 2, 3, and 4 show the least common, somewhat common, and most common words in the dataset, respectively.

### 2.5.1 Selecting a Threshold

Since the provided dataset does not contain stop words, I decided not to filter out common words. The most common words in the dataset appear to contain useful information for clustering.

4

Table 2: 10 least common words with the number and percentage of documents that they appear in.

| Word | Documents | Percentage of Documents |
|---|---|---|
| frazioni | 0 | 0.0% |
| threeletter | 1 | 0.0062881% |
| budjovice | 1 | 0.0062881% |
| gmina | 1 | 0.0062881% |
| ortsgemeinden | 1 | 0.0062881% |
| headwater | 2 | 0.012576% |
| baronetage | 3 | 0.018864% |
| breaststroke | 3 | 0.018864% |
| voronezh | 3 | 0.018864% |
| rosettes | 3 | 0.018864% |

Table 3: 10 words near the median with the number and percentage of documents that they appear in.

| Word | Documents | Percentage of Documents |
|---|---|---|
| multimillion | 16 | 0.10061% |
| pisa | 16 | 0.10061% |
| pranks | 16 | 0.10061% |
| pesticides | 16 | 0.10061% |
| hesitation | 16 | 0.10061% |
| convection | 16 | 0.10061% |
| ortiz | 16 | 0.10061% |
| stagnation | 16 | 0.10061% |
| gonzlez | 16 | 0.10061% |
| cummins | 16 | 0.10061% |

Table 4: 10 most common words with the number and percentage of documents that they appear in.

| Word | Documents | Percentage of Documents |
|---|---|---|
| well | 4009 | 25.209% |
| second | 3538 | 22.247% |
| high | 2836 | 17.833% |
| family | 2579 | 16.217% |
| group | 2412 | 15.167% |
| north | 2364 | 14.865% |
| major | 2298 | 14.45% |
| large | 2227 | 14.004% |
| general | 2187 | 13.752% |
| long | 2164 | 13.607% |

Selecting a lower threshold involves finding the right balance between information loss and cluster sparsity. Words that are only used in one document should be ignored since they provide no useful clustering information. The problem gets murkier with words that are used in just a few datasets. Table 3 lists ten words that are somewhat common, in that their counts are exactly at the median word count. Some of these words seem useful for clustering, such as "multimillion", "pesticides", and "convection". Others are names which do not necessarily indicate document similarity.

One of the large problems here is that we are working with only 15,000 of the 5 million English Wikipedia articles. If we were working with more data, we could produce more fine-grained clusters with rarer words. However, since we only have a tiny subset of Wikipedia, we just don't have enough examples to produce these narrow clusters.

There is no correct way to choose a tf-idf threshold. I ended up choosing to get rid of the least common half of words, but many other choices would be just as valid.

## 2.6 Effect of Seed on Distortion

After observing some unexpected variability in my results, I decided to investigate the effect of random seeds on distortion.

I found that the choice of seed can have a large impact on cluster quality and distortion. Because of this, it is very important to run K-Means with at least a few different seeds and select the best one. I opted to not do this in the above experiments because it would have been too computationally expensive, but I will try many different seeds when generating the final clusters.
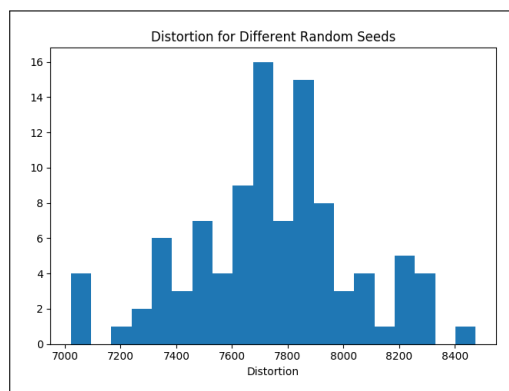


Figure 5: Distortion using K-Means++ on 100 training documents with different random seeds.

# 3 Conclusion

Table 5 shows the final clusters. While 64 clusters were generated, only the largest 10 are shown. These clusters with words that occurred in at least 16 documents, which is about half of all words. K-Means++ was used to generate the initial cluster centroids. Names have been added to each cluster.

For the most part, the clusters make sense.

## 3.1 Issues with K-Means

### 3.1.1 Variability in Cluster Size

The most glaring issue with these results is that the sizes of the clusters vary greatly. One cluster contains about half of the documents while the other 63 contain the other half. There are many useless one document clusters. It is difficult to know if this is intrinsic to the data, or an issue with the K-Means algorithm. Since we are working with over 10,000 dimensions, it is impossible to visualize the data and get an intuition about which model would work best. Using dimensionality

### 3.1.2 Different Meanings of the Same Word

For example, cluster 9 contains documents about mars. Its main words are "mars", "crater", "astronomical", "expedition", and "planets". However, there are two documents in this cluster that have nothing to do with mars. "hostagefilm" (*Hostage* (film)) is an American action movie whose main character is named Marshall "Mars" Krupcheck. "marsrapper" (*Mars* (rapper)) is a Bay Area rapper who goes by "Mars".

### 3.2 Further Research

One of the core assumptions of K-Means is that the data can be grouped into spherically symmetric clusters of the same size. This is quite a strong assumption, and it does not seem to be true in the case of Wikipedia articles. It would be interesting to experiment with other clustering techniques better suited to high dimensional text data.

## 4 Conclusion

Table 5: 10 largest K-Means clusters with $K = 64$. Using words that appear in at least 16 (0.1%) documents with K-Means++ for initializing cluster centroids.

| Cluster | Size | Words | Documents | Label |
|---|---|---|---|---|
| 0 | 8486 | females<br>station<br>family<br>north<br>located | kotavamsa<br>tornadoesof1993<br>whitby<br>colinmeads<br>mladenjiiubiofbribir | |
| 1 | 2226 | students<br>system<br>high<br>program<br>technology | edmondscommunitycollege<br>helderbergcollege<br>stargateschool<br>miltonhighschoolmiltongeorgia<br>govthazimuhammadmohsincollege | Colleges |
| 2 | 2085 | church<br>published<br>daughter<br>royal<br>paris | molire<br>oberlincongregationalchurch<br>lundbyoldchurch<br>stmaryschurchgrodno<br>dioceseoflimerickandkillaloe | Churches |
| 3 | 1001 | party<br>served<br>member<br>general<br>senate | partyidentification<br>labourfarmerparty<br>serbianliberalparty<br>bronwenmaher<br>democraticalliancesouthafrica | Political Parties |
| 4 | 932 | season<br>club<br>playing<br>seasons<br>player | bjsamsamericanfootball<br>dancingwiththestars<br>davidmccracken<br>johnflanaganfootballer<br>gilbertcurgenven | Sports |
| 5 | 640 | album<br>released<br>songs<br>records<br>rock | primalscream<br>thegreatestdaytakethatalbum<br>conflictingemotions<br>bornthisway<br>sweetkisses | Rock Music |
| 6 | 280 | japanese<br>japan<br>chinese<br>pearl<br>characters | frederickringer<br>astorhousehotelshanghai19221959<br>japanesebadger<br>imperialjapanesearmyairforce<br>listofflclepisodes | Imperial Japan |
| 7 | 44 | pop<br>songs<br>album<br>chart<br>rock | popmusic<br>teenpop<br>britishpopmusic<br>talkinginyoursleepcrystalgaylesong<br>blahblahblahalbum | Pop Music |
| 8 | 27 | investigation<br>money<br>system<br>june<br>doctor | digitalmonetarytrust<br>andyhayman<br>martensvillesatanicsexscandal<br>johnlittlechild<br>unitedstatesvlibby | Crime |
| 9 | 22 | mars<br>crater<br>astronomical<br>expedition<br>planets | hostagefilm<br>marsrapper<br>entomopter<br>mars2<br>carbonatesonmars | Mars |