

# Practical Issues of Clustering Relatively Small Text Data Sets for Business Purposes

Nikita Nikitinsky, Tamara Sokolova and Ekaterina Pshehotskaya  
InfoWatch Company

Russia

Nikita.Nikitinsky@infowatch.com, Tamara.Sokolova@infowatch.com,  
Ekaterina.Pshehotskaya@infowatch.com

## ABSTRACT

Clustering of relatively small sets of documents has become a frequent task in small business. Current topic modeling and clustering algorithms can handle this task, but there are some ways to improve the quality of cluster analysis, for example, by introducing some combined algorithms.

In this paper, we will conduct some experiments to define the best clustering algorithm among LSI, LDA and LDA+GS combined with GMM and find heuristics to improve the performance of the best algorithm.

## KEYWORDS

Clustering, cluster analysis, topic modeling, LDA, LSI, GMM, Silhouette Coefficient

## 1 INTRODUCTION

Cluster analysis of relatively small sets of documents (up to 50 000) is a task, which may be essential in small business. It might be necessary to cluster, for example, weekly document stream for DLP (Data Leakage Prevention) purposes (e.g. easier categorization of documents).

To cluster small sets of documents we primarily need high clustering quality and may pay little attention to speed or computational complexity of a clustering algorithm – obviously, because modern computer hardware allows the user to perform complex computations in a short time, so small data sets are clustered fast even if an algorithm with high computational complexity is used. That is why we decided to conduct some experiments on algorithms with high computational complexity in order to combine them in a way, allowing us to maximize quality of clustering.

## 2 METHODS

Cluster analysis or clustering is a convenient method for identifying homogenous groups of objects called clusters. Objects in a specific cluster share many characteristics, but are very dissimilar to objects not belonging to that cluster. Further in this paper we will discuss clustering algorithms where every object can belong only to one cluster - including cases where an object may belong to no cluster at all. In such cases we will create so called «garbage» cluster and put there all objects not classified by an algorithm.

We will use the following clustering and topic modeling algorithms to create a combination showing highest performance:

LSI (Latent Semantic Indexing) — is an unsupervised machine learning method, which is mostly used for dimensionality reduction. It is an indexing and retrieval method that uses a mathematical technique called singular value decomposition (SVD) to identify patterns in the relationships between the terms and concepts contained in an unstructured collection of text. LSI is based on the principle that words that are used in the same contexts tend to have similar meanings. A key feature of LSI is its ability to extract the conceptual content of a body of text by establishing associations between those terms that occur in similar contexts. [1]

LDA (Latent Dirichlet Allocation) - is also an unsupervised machine learning method, which is mostly used for object clustering. It is a generative model that allows sets of observations to be explained by unobserved

groups that explain why some parts of the data are similar. For example, if observations are words collected into documents, it posits that each document is a mixture of a small number of topics and that each word's creation is attributable to one of the document's topics. Computational complexity of LDA+GS is  $O(NKW)$  where  $N$  is a number of documents,  $K$  is a number of clusters and  $W$  is the number of words in vocabulary. [2]

Although mentioned above methods can be used alone, we will conduct experiments, in which we combine them with the following algorithms:

**GMM Classifier** (Gaussian Mixture Model), which is an unsupervised machine learning method, is a probabilistic model that assumes all the data points are generated from a mixture of a finite number of Gaussian distributions with unknown parameters. One can think of mixture models as generalizing k-means clustering to incorporate information about the covariance structure of the data as well as the centers of the latent Gaussians. Computational complexity of GMM (using EM-algorithm for convergence) is  $O(tkmn^3)$  where  $k$  is the number of clusters,  $n$  is the number of dimensions in a sample,  $m$  is a number of samples and  $t$  is a number of iterations. [3]

We decided to select GMM as the most appropriate for our experiments because it is considered a versatile modeling tool for cluster analysis and its performance is much higher compared to, for example, K-means.

When applying GMM we arrange every object only to one cluster (thus, we make it easier to estimate overall performance).

**GS** (Gibbs Sampling) is a Markov chain Monte Carlo (MCMC) algorithm for obtaining a sequence of observations which are approximated from a specified multivariate probability distribution, when direct sampling is difficult. GS is widely used to enhance quality of topic modeling algorithms; it is a good algorithm for processing when the dimension of data is not very high. With high dimensional data it may be better to use Variational EM algorithm. [4] In our experiments we applied faster version

of GS algorithm named Collapsed Gibbs Sampling algorithm.

### 3 EVALUATION METRICS

To evaluate algorithm performance we used two types of metrics often utilized for cluster analysis purposes:

#### 3.1 External Evaluation Metrics

In external evaluation, clustering results are evaluated based on data that was not used for clustering, such as known class labels and external benchmarks. Such benchmarks consist of a set of pre-classified items, and these sets are often created by human (experts). Thus, the benchmark sets can be thought of as a gold standard for evaluation. [5]

We used the following external measurements:

**Jaccard index** - also known as the Jaccard similarity coefficient, is a statistic used for comparing the similarity and diversity of sample sets. The Jaccard coefficient measures similarity between finite sample sets, and is defined as the size of the intersection divided by the size of the union of the sample sets. [6]

**V-measure score** - is an entropy-based measure which explicitly measures how successfully the criteria of homogeneity and completeness have been satisfied. V-measure is computed as the harmonic mean of distinct homogeneity and completeness scores, just as precision and recall are commonly combined into F-measure. [7]

**Adjusted Rand score** - is a measure of the similarity between two data clusterings.[8]

**Adjusted mutual information score** - a variation of mutual information (which is a measure of the variables' mutual dependence) may be used for comparing clusterings. [9][10]

### 3.2 Internal Evaluation Metrics

In internal evaluation clustering result is evaluated based on the data that was clustered itself. These methods usually assign the best score to the algorithm that produces clusters with high similarity within a cluster and low similarity between clusters. [11]

We used the following internal measurement:

**Silhouette Coefficient** — is a measure of how appropriately the data has been clustered and how well each object lies within its cluster.

The Silhouette Coefficient is defined for each sample and is composed of two scores:

1. The mean distance between a sample and all other points in the same class.
2. The mean distance between a sample and all other points in the next nearest cluster.

We used cosine metric as the most common for measuring the distances for Silhouette Coefficient. When we have higher value of Silhouette Coefficient, it means that we have better distribution of documents to topics. [12]

Based on Silhouette Coefficient measurements we apply Elbow method to define the number of clusters. This method assumes a choice of a number of clusters so that adding another cluster doesn't give much better modeling of the data (so called “Knee of a curve”). This method was originally designed to make predictions based on the percentage of variance explained and in some cases may appear unsuitable; in such cases we will choose the number of clusters where Silhouette Coefficient reaches maximum value [13]

Since, in real conditions, we are unable to use external metrics for evaluation of algorithms (because we usually don't know the true number of clusters), we will evaluate quality of our models basing mostly on Silhouette Coefficient, applying external metrics as supplementary.

## 4 DATA SETS

We used some different data sets to check and validate the results:

1. Data set containing 600 documents, distributed to 5 topics – a «good» collection (distribution of documents: 83 to 163 documents per topic). Topics are easily distinguishable by human expert.
2. Data set containing 157 documents, distributed to 14 topics - «bad» collection (distribution of documents: 3 to 21 documents per topic). Topics are not distinguishable by human expert.
3. Data set containing 1000 documents, randomly assigned from the real document stream of the company; topic distribution is not predetermined; human experts considered the number of topics between 3 and 5 (including 3 and 5).
4. Data set containing 35000 documents, randomly assigned from the real document stream of the company; topic distribution is not predetermined. Human experts then estimated quality of the best algorithm performance on this data set.

## 5 EXPERIMENTS

We tested all these algorithms on the «good» collection to find out the best one and then evaluated the best algorithm performance on other collections

### 5.1 Choosing the Best Algorithm

#### 5.1.1 LSI+GMM

Data preprocessing:

All words with length less than 3 symbols were deleted as well as all non-alphabetic characters.

To obtain better results we preprocessed input data with TF-IDF algorithm.

In this algorithm we may vary two main parameters: number of LSI topics and number of GMM clusters.

The LSI algorithm takes as input the

collection of documents, processes it and then documents-topic matrix is returned. This matrix is then given to an input of GMM classifier, which processes the input matrix assembling documents to final categories (this is likely to increase the quality of clustering).

We tested two heuristics:

1. Number of LSI topics is equal to number of output GMM clusters
2. Number of LSI topics is equal to number of output GMM clusters plus one, such as number of LSI topics is  $n+1$ , while number of GMM clusters is  $n$  (one of the topics becomes so called «garbage» topic — it accumulates objects, which could not be unambiguously arranged to other «real» topics)

Table 1 contains evaluation metrics estimated on the «good» collection for LSI+GMM algorithm with 5 output categories:

**Table 1.**

	Heuristic 1	Heuristic 2
Jaccard index	0.575	0.57
Adjusted mutual information score	0.75	0.735
Adjusted Rand score	0.66	0.66
V measure score	0.74	0.74
Silhouette Coefficient	0.61	0.5

We can see that both heuristics showed comparable results when tested on a real number of categories; Heuristic 2 showed a decrease in Silhouette Coefficient value.

But, more generally, if we vary the number of output categories and estimate Silhouette Coefficient for them we will get the following results (Figures 1, 2):



**Figure 1. LSI+GMM, Heuristic 1**



**Figure 2. LSI+GMM, Heuristic 2**

According to the results, the Silhouette Coefficient reached higher levels when we implemented Heuristic 2 (Figure 2). Nevertheless, both pikes indicated incorrect number of output clusters (6 and 8 correspondingly).

### 5.1.2. LDA+GMM

Data preprocessing:

All words with length less than 3 symbols were deleted as well as all non-alphabetic characters

Words occurring only once (hapax legomena) were deleted

The LDA algorithm takes as input the collection of documents, processes it and then

documents-topic matrix is returned. This matrix is then given to an input of GMM classifier which processes the input matrix assembling documents to final clusters (this must increase the quality of clustering).

We tested the same two heuristics.

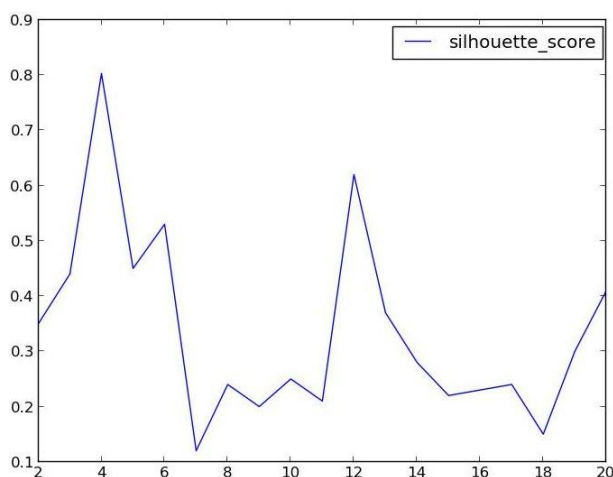
Table 2 contains metrics estimated on the «good» collection for LDA+GMM algorithm with 5 output categories:

**Table 2.**

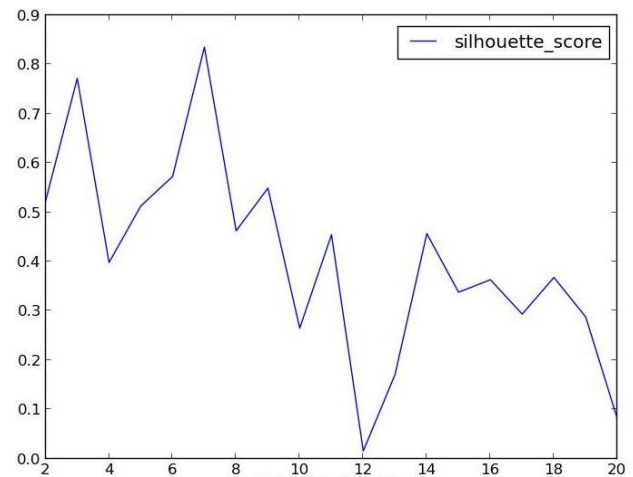
	Heuristic 1	Heuristic 2
Jaccard index	0.51	0.85
Adjusted mutual information score	0.57	0.83
Adjusted Rand score	0.53	0.76
V measure score	0.6	0.84
Silhouette Coefficient	0.45	0.52

We can see that Heuristic 2 showed far better results for external metrics, but insignificantly better result for Silhouette Coefficient.

If we vary the number of output categories and estimate Silhouette Coefficient for them we will get the following results (Figures 3, 4):



**Figure 3. LDA+GMM, Heuristic 1**



**Figure 4. LDA+GMM, Heuristic 2**

According to the results, Silhouette Coefficient reached a bit higher levels when we implemented Heuristic 2 (Figure 4). Nevertheless, both pikes indicated incorrect number of output clusters (4 and 7 correspondingly).

### 5.1.3 LDA+GS+GMM

Data preprocessing:

All words with length less than 3 symbols were deleted as well as all non-alphabetic characters

Words occurring only once (hapax legomena) were deleted

In this algorithm we may vary three main parameters: number of LDA topics, number of Gibbs Samples and number of GMM clusters.

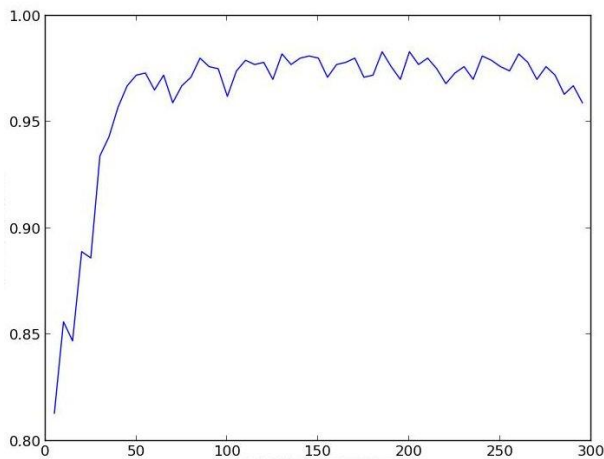
For given quantity of LDA topics there are n iterations of Gibbs Sampling (where n is number of Gibbs Samples) and then documents-topic matrix is returned. This matrix is then given to an input of GMM classifier which processes the input matrix assembling documents to final clusters.

Choosing proper number of Gibbs Samples:

Knowing the real quantity of output categories we iteratively start the algorithm changing the number of samples and keeping other parameters the same.

The best number of Gibbs Samples is considered the number of samples when metric (e.g, Silhouette Coefficient) reaches

highest values and then doesn't fluctuate much.



**Figure 5.**

We selected the best number of GS samples on the “good” collection. The unchanged parameters were the number of LDA topics and the number of GMM clusters (as in Heuristic 2). As we can see from the picture (Figure 5), the plotted line reaches highest values at 50 samples and then don't fluctuate much, so we can choose any quantity of samples above 50, so we will then use 100 samples as optimal and versatile number of samples.

We tested the same two heuristics.

Table 3 contains metrics estimated on the «good» collection for LDA+GS+GMM algorithm with 5 output categories:

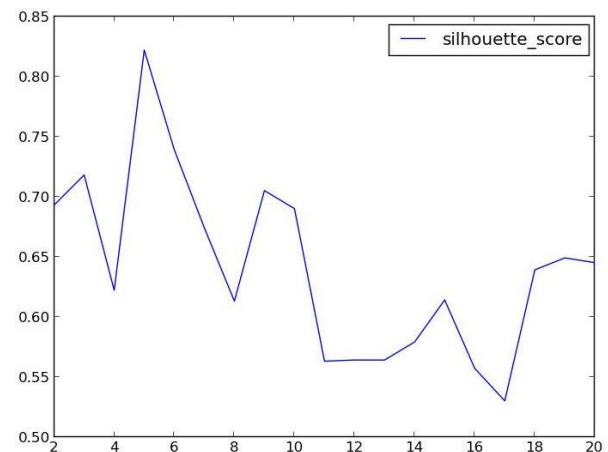
**Table 3.**

	Heuristic 1	Heuristic 2
Jaccard index	0.66	0.99
Adjusted mutual information score	0.77	0.99
Adjusted Rand score	0.72	0.99
V measure score	0.79	0.99
Silhouette Coefficient	0.82	0.98

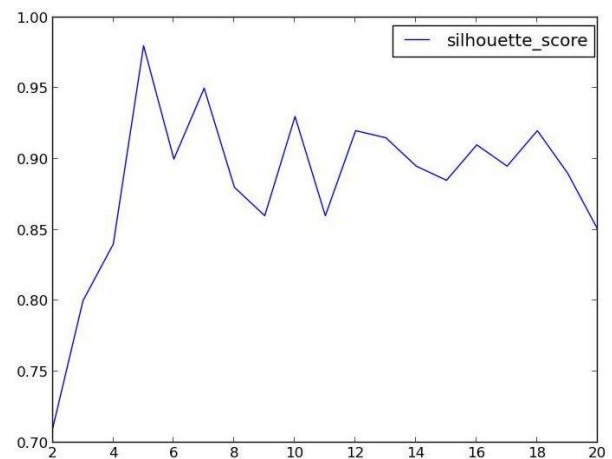
We can see that Heuristic 2 showed far better results for all metrics. It means that documents are better distributed to said number of output categories with Heuristic 2 implemented for this algorithm.

If we vary the number of output categories and estimate Silhouette Coefficient for them

we will get the following results (Figures 6, 7):



**Figure 6. LDA+GS+GMM, Heuristic 1**



**Figure 7. LDA+GS+GMM, Heuristic 2**

According to the results, while both pikes indicated the same true number of clusters, Silhouette Coefficient reached higher levels when we implemented Heuristic 2 (Figure 7). We can suggest that Heuristic 2 improves the performance of LDA+GS+GMM and intensifies the results making it easier to determine the number of output categories

## 5.2 Estimating the Best Algorithm on Other Data Sets

We tested LDA+GS+GMM algorithm on other collections using the parameters that we considered the best testing the algorithm on the “good” collection:

Number of GS samples is equal to 100

Number of LDA topics is equal to number of GMM clusters plus one (e.g. while number of GMM clusters is 5, number of LDA topics is 6)



Data preprocessing:

All words with length less than 3 symbols were deleted as well as all non-alphabetic characters

Words occurring only once (hapax legomena) were deleted

### 5.2.1 Data Set №2

We tested LDA+GS+GMM algorithm on the «bad» collection and had the following results (Figure 8):

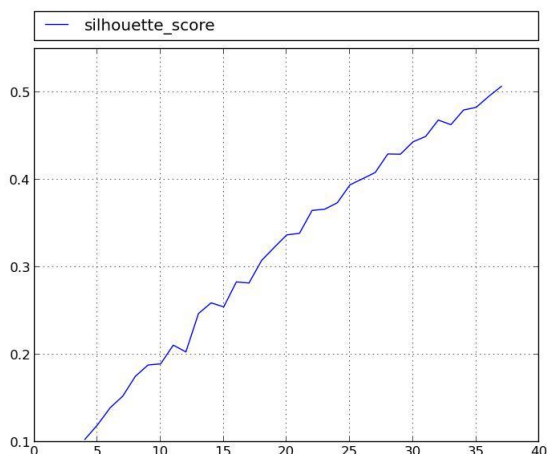


Figure 8.

Assuming that we selected the optimal parameters and using Elbow method based on Silhouette Coefficient plot we found it impossible to define (even approximately) the best number of output categories, because:

1. The distribution of documents to topics is conventional (in such cases there are either no much difference in vocabulary between documents of different categories or difference between all documents is too high to group at least some of them into one definite cluster)

2. Number of documents is small

### 5.2.2. Data Set №3

We tested LDA+GS+GMM algorithm on the data set №3 containing 1000 documents and had the following results (Figure 9):

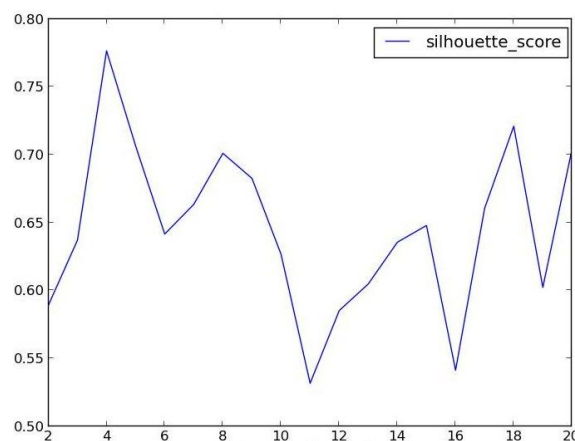


Figure 9.

Basing on Silhouette Coefficient plot we decided that 4 categories is the best number of clusters for this data set. Human experts considered the result of the algorithm good. Documents in four categories could easily be defined as contracts, financial documents, application forms and information letters + instructions.

### 5.2.3. Data Set №4

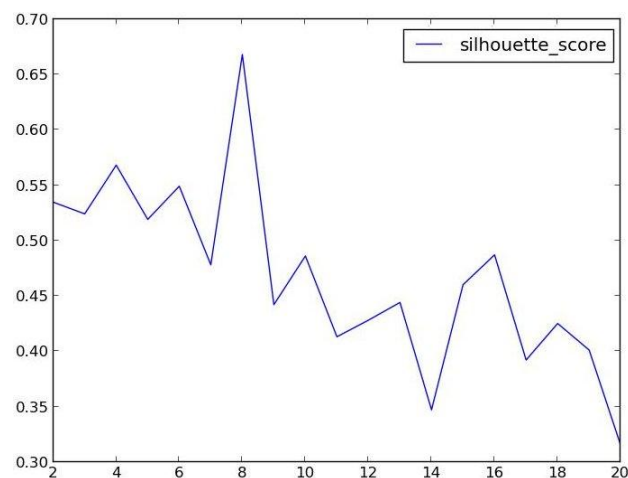


Figure 10.

Basing on Silhouette Coefficient plot (Figure 10) we decided that 8 categories were the best quantity for this data set.

Human experts defined documents in 8 categories as contracts, financial documents, documents in other languages, information letters, instructions, application forms and other internal documents.

## 6 CONCLUSION

According to the experiments we conducted, the best algorithm for processing relatively small set of documents (up to 50 000) with

relatively small quantity of topics (up to 20) is LDA+GS+GMM. The Heuristic 2 may help to improve quality of LDA+GS+GMM and make it easier to determine number of output categories. Usage of Silhouette Coefficient is considered appropriate for determining best number of output clusters. The data set should not be too small in order to provide the clustering algorithm with processable data: data sets containing less than 500 documents are likely to be incorrectly classified.

## 7 FURTHER READING

There are some papers on automated number of clusters detection algorithms, such as [14], proposing state-of-the-art algorithms that may be useful for cluster analysis.

Although Latent Dirichlet Allocation works well for topic modeling there are now conducted multiple researches on more advanced topic modeling algorithms such as Higher-order Latent Dirichlet Allocation and other Higher-order topic modeling algorithms [15].

## 8 REFERENCES

- [1] Deerwester, S., Dumais S., Landauer T., Furnas G., Beck L., Improving Information Retrieval with Latent Semantic Indexing, Proceedings of the 51st Annual Meeting of the American Society for Information Science 25, 1988, pp. 36–40.
- [2] Blei, D. M.; Ng, A. Y.; Jordan, M. I. (January 2003). "Latent Dirichlet allocation". In Lafferty, John. *Journal of Machine Learning Research* 3 (4–5): pp. 993–1022
- [3] Bishop C.M. (2006) *Pattern recognition and machine learning*. Springer, Berlin
- [4] Casella, G., Edward, I. (1992). "Explaining the Gibbs sampler". *The American Statistician* 46 (3): 167–174
- [5] Kaufman L, Rousseeuw P.J. (2005) *Finding groups in data. An introduction to cluster analysis*. Wiley, Hoboken, NY
- [6] Tan, P.-N.; Steinbach, M.; Kumar, V. (2005), *Introduction to Data Mining*
- [7] Rosenberg, A. and Hirschberg, J.. "V-Measure: A conditional entropy-based external cluster evaluation measure". *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*: 410–420
- [8] Rand, W. M. (1971). "Objective criteria for the evaluation of clustering methods". *Journal of the American Statistical Association (American Statistical Association)* 66 (336): 846–850
- [9] Meila, M. (2007). "Comparing clusterings—an information based distance". *Journal of Multivariate Analysis* 98 (5): 873–895.
- [10] Vinh, N. X.; Epps, J.; Bailey, J. (2009). "Information theoretic measures for clusterings comparison". *Proceedings of the 26th Annual International Conference on Machine Learning - ICML 09*. p. 1.
- [11] Manning, C. D, Raghavan, P. & Schütze, H., *Introduction to Information Retrieval*. Cambridge University Press.
- [12] Rousseeuw, P, J. (1987). "Silhouettes: a Graphical Aid to the Interpretation and Validation of Cluster Analysis". *Journal of Computational and Applied Mathematics* 20: 53–65
- [13] Ketchen, D, J., Jr & Shook, C, L. (1996). "The application of cluster analysis in Strategic Management Research: An analysis and critique". *Strategic Management Journal* 17 (6): 441–458.
- [14] Salvador, S, and Chan, P., Determining the Number of Clusters/Segments in Hierarchical Clustering/Segmentation Algorithms, *Proceedings of 16<sup>th</sup> IEEE International Conference on Tools with AI*. Washington, DC: IEEE Computer Society.
- [15] Nelson, C., Pottenger, W. M., Keiler, H., and Grinberg, N., "Nuclear Detection Using Higher-Order Topic Modeling." *2012 IEEE International Conference on Technologies for Homeland Security*. Waltham, MA. 13-15 Nov 2012.