# Text Clustering

Data preparation and pre-processing

Data Transformation

Text Clustering

Model Evaluation

Error Analysis

Need help

# Data Pre-Processing

⚙️ **Removing Stop Words**

✓ **Converting words to lower case**

✂️ **Performing Lemmatization**

📄 **Partition Every book to 200 Partitions**

📜 **Every Partition have 150 words**

# Data Transformation

Bag Of Words(BOW)

TF-IDF

LDA

Word2Vec

# Bag of Words(BOW) Transformation

A bag of words is a representation of text that describes the occurrence of words within a document.

| | aaron | abandon | abandoned | abandoning | abandonment | abated | abb | abbe | abbey | abbott | ... | zodiacal | zonal | zone | zool | zoologique | zoologist | zoology | zur | zwischen | zygomatic |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 995 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 996 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 997 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 998 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 999 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

1000 rows × 16048 columns

# TF-IDF

Term frequency (TF) vectors show how important words are to documents. They are computed by using:

$$tf(term, document) = \frac{number\ of\ times\ the\ term\ occurs\ in\ the\ document}{total\ number\ of\ terms\ in\ the\ document}$$

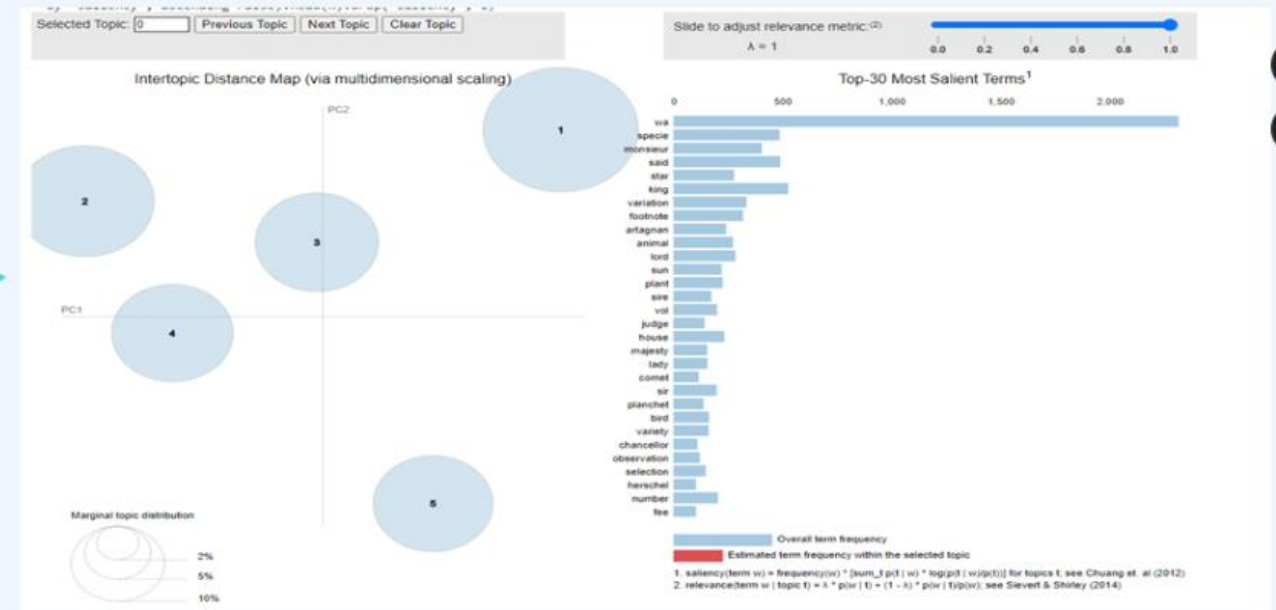| | aaron | abandon | abandoned | abandoning | abandonment | abated | abb | abbe | abbey | abbott | ... | zodiacal | zonal | zone | zool | zoologique | zoologist | zoology | zur | zwischen | zygomatic |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 2 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 3 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 4 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 995 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 996 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 997 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 998 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 999 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

000 rows × 16048 columns

# LDA

LDA is used to classify text in a document to a particular topic. It builds a topic per document model and words per topic model, modeled as Dirichlet distributions. Each document is modeled as a multinomial distribution of topics and each topic is modeled as a multinomial distribution of words.

| | 1 | 2 | 3 | 4 | 5 | res |
|---|---|---|---|---|---|---|
| **0** | 40.195923 | 0.179589 | 86.287994 | 17.819510 | 6.548666 | 3 |
| **1** | 6.840732 | 0.179266 | 0.174638 | 9.724366 | 134.112686 | 5 |
| **2** | 42.947739 | 86.382599 | 3.532179 | 5.912141 | 12.257037 | 2 |
| **3** | 13.155953 | 1.853795 | 0.175965 | 131.157608 | 4.688347 | 4 |
| **4** | 11.478949 | 5.613846 | 0.174887 | 15.622274 | 118.141701 | 5 |

```
(array([[ 40.195923  ,   0.17958926,  86.287994  ,  17.81951   ,
          6.548666  ],
        [  6.840732  ,   0.1792659 ,   0.17463806,   9.724366  ,
        134.11269   ],
        [ 42.94774   ,  86.3826    ,   3.5321789 ,   5.9121413 ,
         12.257037  ],
        ...,
        [ 35.19193   , 115.3203    ,   0.17444904,   0.20968111,
          0.13532138],
        [ 20.348946  , 113.14746   ,   7.5316253 ,   0.20784229,
          9.795807  ],
        [  0.37172854,   0.17994398,   4.8293304 ,  26.066813  ,
        119.58385   ]], dtype=float32), None)
```

**LDA as a Topic Modeling**

# LDA as a Topic Modeling

LDA is an example of topic model and is used to classify text in a document to a particular topic. It builds a topic per document model and words per topic model, modeled as Dirichlet distributions.

# Word2Vec

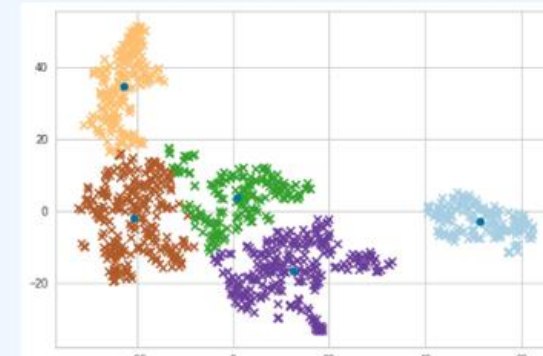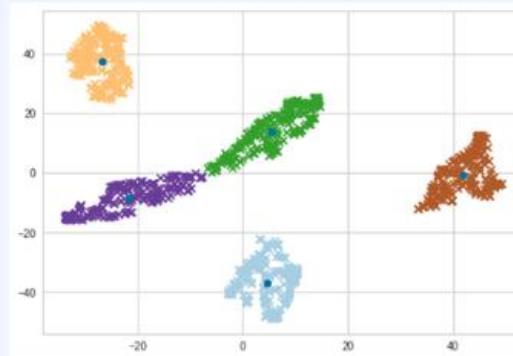Word2Vec consists of models for generating word embedding.
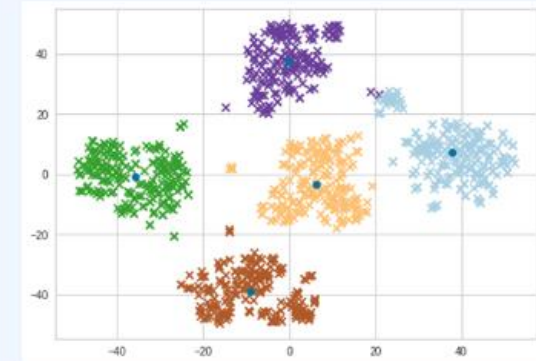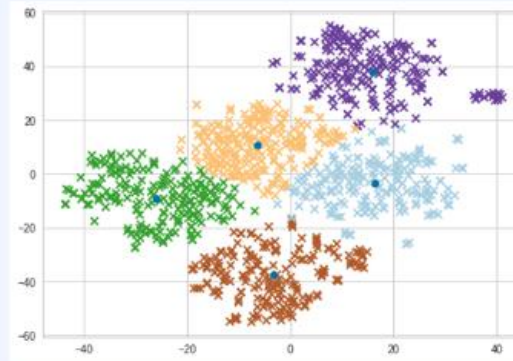
# Text Clustering Algorithms

K-Means Clustering
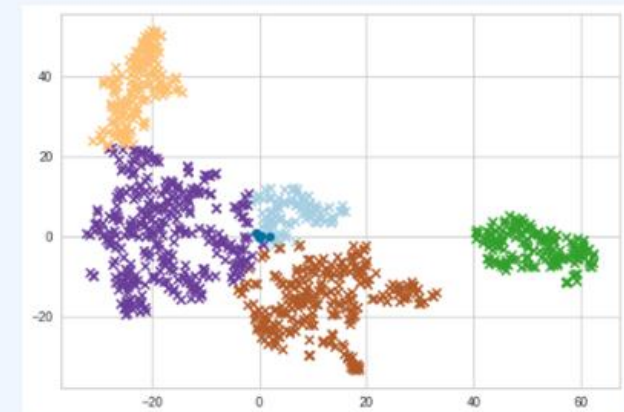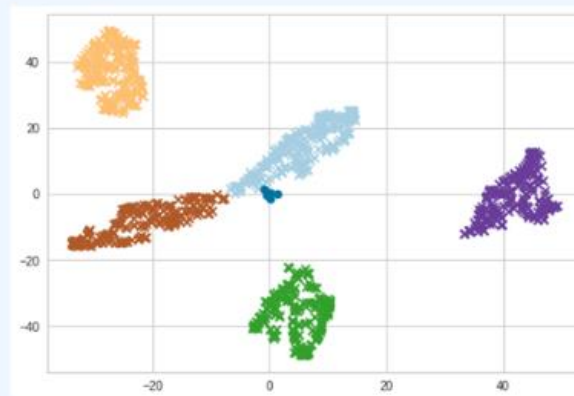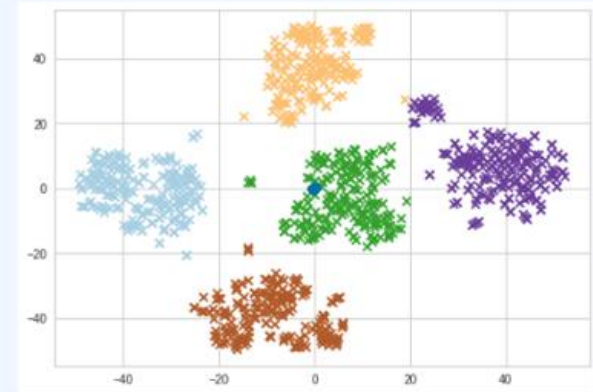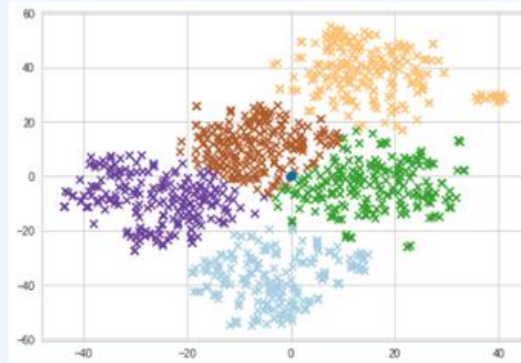
Expectation Maximization

Hierarchical Clustering

# K-Means Clustering

K-Means Clustering is an unsupervised learning algorithm that is used to solve the clustering problems in machine learning or data science. In this topic, we will learn what is K-means clustering algorithm, how the algorithm works, along with the Python implementation of k-means clustering.

# Expectation Maximization

Expectation-Maximization algorithm can be used for the latent variables (variables that are not directly observable and are actually inferred from the values of the other observed variables) too in order to predict their values with the condition that the general form of probability distribution governing those latent variables is known to us. This algorithm is actually at the base of many unsupervised clustering algorithms in the field of machine learning.

# Hierarchical Clustering

Hierarchical clustering, also known as hierarchical cluster analysis, is an algorithm that groups similar objects into groups called clusters. The endpoint is a set of clusters, where each cluster is distinct from each other cluster, and the objects within each cluster are broadly similar to each other.

# Model Evaluation

Model Evaluation using Kappa

Model Evaluation using consistency with V-Score

Model Evaluation using Coherence

Model Evaluation using Silhouette Score

# Model Evaluation using Kappa

Cohen's kappa: a statistic that measures inter-annotator agreement.

| Kappa's Score with K-Means clustering algorithm | | | |
|---|---|---|---|
| K-means with BOW | K-means with TF-IDF | K-means with LDA | K-means with Word2Vec |
| 0.92125 | 0.96375 | 0.98625 | 0.73375 |
| Kappa Score of LDA as Topic Modeling | 0.27795 | | |
| Kappa's Score with Expectation Maximization(EM) algorithm clustering algorithm | | | |
| EM with BOW | EM with TF-IDF | EM with LDA | EM with Word2Vec |
| 0.92750 | 0.96250 | 0.98625 | 0.70125 |
| Kappa's Score with Hierarchal clustering algorithm | | | |
| Hierarchal clustering with BOW | Hierarchal clustering with TF-IDF | Hierarchal clustering with LDA | Hierarchal clustering with Word2Vec |
| 0.91000 | 0.96125 | 0.98625 | 0.72875 |

# Model Evaluation using consistency with V-Score

V-measure cluster labeling given a ground truth.

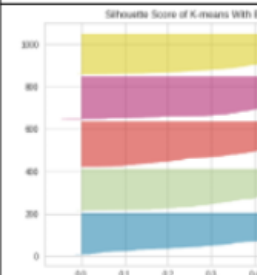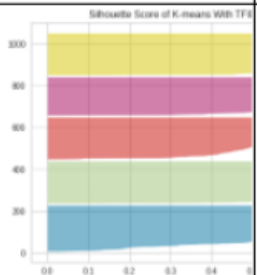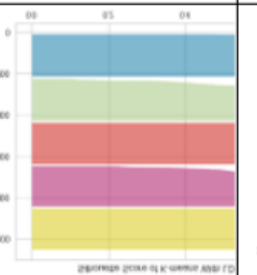| V-Score with K-Means clustering algorithm | | | |
|---|---|---|---|
| **K-means with BOW** | **K-means with TF-IDF** | **K-means with LDA** | **K-means with Word2Vec** |
| 0.84667087642090 49 | 0.92272558192760 63 | 0.96041867936215 67 | 0.60410870337612 3 |
| V-Score with Expectation Maximization(EM) algorithm clustering algorithm | | | |
| **EM with BOW** | **EM with TF-IDF** | **EM with LDA** | **EM with Word2Vec** |
| 0.85236487417651 33 | 0.92124908692837 59 | 0.96041867936215 67 | 0.61686102073448 51 |
| V-Score with Hierarchal clustering algorithm | | | |
| **Hierarchal clustering with BOW** | **Hierarchal clustering with TF-IDF** | **Hierarchal clustering with LDA** | **Hierarchal clustering with Word2Vec** |
| 0.84601454866905 73 | 0.91978072721278 81 | 0.96041867936215 67 | 0.62704087646871 4 |

# Model Evaluation using Coherence

coherence is used to measure how well the topics are extracted.
Coherence score With LDA using c_v: 0.4363
Higher value is better.
Coherence score With LDA using u_mass: -8.0567
Lower value is better.
t that describes the occurrence of words within a document.

# Model Evaluation using Silhouette Score

Silhouette score is used to evaluate the quality of clusters created using clustering algorithms such as K-Means in terms of how well samples are clustered with other samples that are similar to each other. The Silhouette score is calculated for each sample of different clusters. Silhouette Coefficient or silhouette score is a metric used to calculate the goodness of a clustering technique. Its value ranges from -1 to 1.

| Silhouette Score with K-Means clustering algorithm | | | |
|---|---|---|---|
| K-means with BOW | K-means with TF-IDF | K-means with LDA | K-means with Word2Vec |
| 0.47652 | 0.62407 | 0.72321 | 0.51574 |



| Silhouette Score with Expectation Maximization(EM) algorithm clustering algorithm | | | |
|---|---|---|---|
| EM with BOW | EM with TF-IDF | EM with LDA | EM with Word2Vec |
| 0.47449 | 0.62412 | 0.72321 | 0.48069 |

| Silhouette Score with Hierarchal clustering algorithm | | | |
|---|---|---|---|
| Hierarchal clustering with BOW | Hierarchal clustering with TF-IDF | Hierarchal clustering with LDA | Hierarchal clustering with Word2Vec |
| 0.46485 | 0.62326 | 0.72321 | 0.46807 |

# Error Analysis

## The Model Misclassified 64 rows

The Most common collocations with their repeat count in all records that were labelled uncorrectly according to the human label

[(('cost', 'almost'), 6), (('anyone', 'anywhere'), 6), (('english', 'character'), 6), (('character', 'set'), 6), (('almost', 'restriction'), 6), (('distributed', 'proofreading'), 3), (('copy', 'give'), 3), (('anywhere', 'cost'), 3), (('imp', 'tersbourg'), 1), (('louise', 'valliere'), 1), (('license', 'included'), 1), (('language', 'english'), 1), (('encoding', 'iso'), 1), (('give', 'away'), 1), (('de', 'sci'), 1), (('clearly', 'reader'), 1), (('extent', 'independently'), 1), (('date', 'january'), 1), (('excess', 'defect'), 1), (('heart', 'mind'), 1)]

| | PartitionsList | Label_of_Book | index | clustersOutput |
|---|---|---|---|---|
| 177 | vol rosenberger calculated though lived lynn o... | e | 4 | 0 |
| 3 | may aware project gutenberg ha involved writin... | d | 3 | 0 |
| 189 | dim perception already arrived perhaps observa... | a | 0 | 4 |
| 71 | body much exception part agree large marked in... | c | 2 | 4 |
| 81 | considerable number spread varying distance si... | c | 2 | 4 |
| ... | ... | ... | ... | ... |
| 174 | beneficial le beneficial le size body would be... | c | 2 | 4 |
| 25 | wa thus recognised domain far reaching specula... | e | 4 | 1 |
| 50 | trans vol footnote ibid vol cvii footnote bull... | e | 4 | 0 |
| 175 | dissipation extinction footnote footnote allge... | e | 4 | 0 |
| 115 | undoubtedly wa accads clear authentic insight ... | a | 0 | 2 |

64 rows × 4 columns

**The most common words that threw the machine off**

Prezi

# The most common words that threw the machine off



Common Bigrams in instances that threw the machine off