

Master the Art of Data Science: A Step-by-Step Guide to Becoming a Full Stack Data Scientist

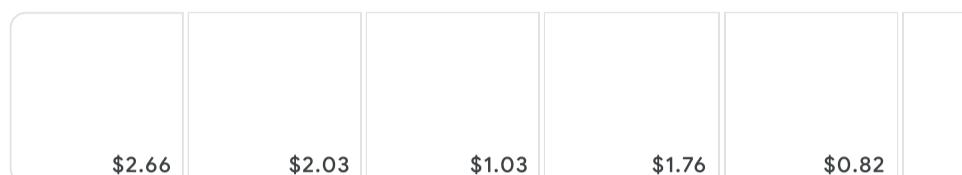
[Explore Now](#)

[Home](#)

# Beginners Guide to Topic Modeling in Python

Shivam5992 Bansal – Published On August 24, 2016 and Last Modified On December 23rd, 2020

[Data Science](#) [Intermediate](#) [NLP](#) [Python](#) [Technique](#) [Text](#) [Topic Modeling](#) [Unstructured Data](#) [Unsupervised](#)



## Introduction

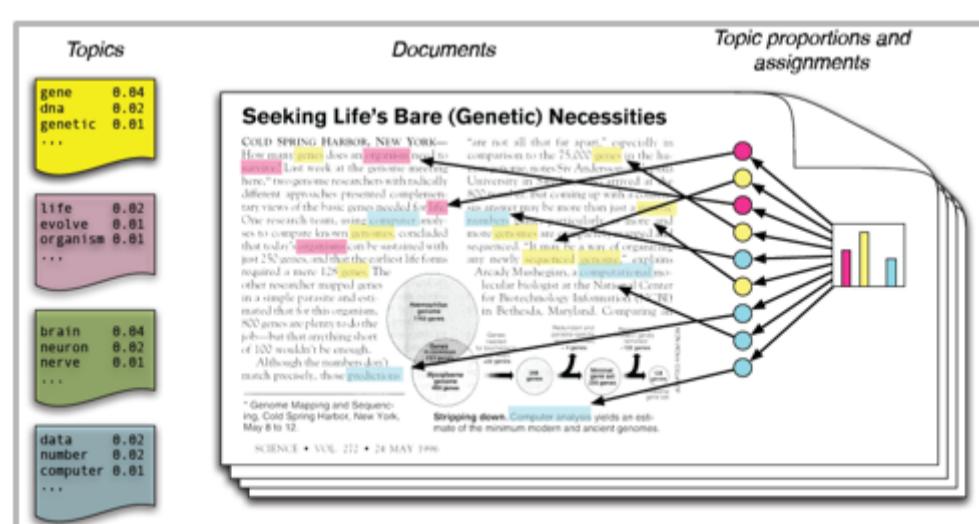
Analytics Industry is all about obtaining the “Information” from the data. With the growing amount of data in recent years, that too mostly unstructured, it’s difficult to obtain the relevant and desired information. But, technology has developed some powerful methods which can be used to mine through the data and fetch the information that we are looking for.

One such technique in the field of text mining is Topic Modelling. As the name suggests, it is a process to automatically identify topics present in a text object and to derive hidden patterns exhibited by a text corpus. Thus, assisting better decision making.

Topic Modelling is different from rule-based text mining approaches that use regular expressions or dictionary based keyword searching techniques. It is an unsupervised approach used for finding and observing the bunch of words (called “topics”) in large clusters of texts.

Topics can be defined as “a repeating pattern of co-occurring terms in a corpus”. A good topic model should result in – “health”, “doctor”, “patient”, “hospital” for a topic – Healthcare, and “farm”, “crops”, “wheat” for a topic – “Farming”.

Topic Models are very useful for the purpose for document clustering, organizing large blocks of textual data, information retrieval from unstructured text and feature selection. For Example – New York Times are using topic models to boost their user – article recommendation engines. Various professionals are using topic models for recruitment industries where they aim to extract latent features of job descriptions and map them to right candidates. They are being used to organize large datasets of emails, customer reviews, and user social media profiles.



So, if you aren't sure about the complete process of topic modeling, this guide would introduce you to various concepts followed by its implementation in python.

## Table of Content

- Latent Dirichlet Allocation for Topic Modeling
  - Parameters of LDA
- Python Implementation
  - Preparing documents
  - Cleaning and Preprocessing
  - Preparing document term matrix
  - Running LDA model
  - Results
- Tips to improve results of topic modelling
  - Frequency Filter
  - Part of Speech Tag Filter
  - Batch Wise LDA
- Topic Modeling for Feature Selection

## Latent Dirichlet Allocation for Topic Modeling

There are many approaches for obtaining topics from a text such as – Term Frequency and Inverse Document Frequency. NonNegative Matrix Factorization techniques. Latent Dirichlet Allocation is the most popular topic modeling technique and in this article, we will discuss the same.

LDA assumes documents are produced from a mixture of topics. Those topics then generate words based on their probability distribution. Given a dataset of documents, LDA backtracks and tries to figure out what topics would create those documents in the first place.

LDA is a matrix factorization technique. In vector space, any corpus (collection of documents) can be represented as a document-term matrix. The following matrix shows a corpus of N documents D1, D2, D3 ... Dn and vocabulary size of M words W1,W2 .. Wn. The value of i,j cell gives the frequency count of word Wj in Document Di.

	W1	W2	W3	Wn
D1	0	2	1	3
D2	1	4	0	0
D3	0	2	3	1
Dn	1	1	3	0

LDA converts this Document-Term Matrix into two lower dimensional matrices – M1 and M2.

M1 is a document-topics matrix and M2 is a topic – terms matrix with dimensions (N, K) and (K, M) respectively, where N is the number of documents, K is the number of topics and M is the vocabulary size.

	K1	K2	K3	K
D1	1	0	0	1
D2	1	1	0	0
D3	1	0	0	1
Dn	1	0	1	0

	W1	W2	W3	Wm
K1	0	1	1	1
K2	1	1	1	0
K3	1	0	0	1
K	1	1	0	0

Notice that these two matrices already provides topic word and document topic distributions. However, these distribution needs to be improved, which is the main aim of LDA. LDA makes use of sampling techniques in order to improve these matrices.

It Iterates through each word "w" for each document "d" and tries to adjust the current topic – word assignment with a new assignment. A new topic "k" is assigned to word "w" with a probability P which is a product of two probabilities p1 and p2.

For every topic, two probabilities p1 and p2 are calculated.  $P_1 = p(\text{topic } t / \text{document } d)$  = the proportion of words in document d that are currently assigned to topic t.  $P_2 = p(\text{word } w / \text{topic } t)$  = the proportion of assignments to topic t over all documents that come from this word w.

The current topic – word assignment is updated with a new topic with the probability, product of p1 and p2 . In this step, the model assumes that all the existing word – topic assignments except the current word are correct. This is essentially the probability that topic t generated word w, so it makes sense to adjust the current word's topic with new probability.

After a number of iterations, a steady state is achieved where the document topic and topic term distributions are fairly good. This is the convergence point of LDA.

## Parameters of LDA

Alpha and Beta Hyperparameters – alpha represents document-topic density and Beta represents topic-word density. Higher the value of alpha, documents are composed of more topics and lower the value of alpha, documents contain fewer topics. On the other hand, higher the beta, topics are composed of a large number of words in the corpus, and with the lower value of beta, they are composed of few words.

Number of Topics – Number of topics to be extracted from the corpus. Researchers have developed approaches to obtain an optimal number of topics by using Kullback Leibler Divergence Score. I will not discuss this in detail, as it is too mathematical. For understanding, one can refer to this<sup>[1]</sup> original paper on the use of KL divergence.

Number of Topic Terms – Number of terms composed in a single topic. It is generally decided according to the requirement. If the problem statement talks about extracting themes or concepts, it is recommended to choose a higher number, if problem statement talks about extracting features or terms, a low number is recommended.

Number of Iterations / passes – Maximum number of iterations allowed to LDA algorithm for convergence.

[You can learn topic modeling in depth here.](#)

## Running in python

### Preparing Documents

Here are the sample documents combining together to form a corpus.

```
doc1 = "Sugar is bad to consume. My sister likes to have sugar, but not my father."
doc2 = "My father spends a lot of time driving my sister around to dance practice."
doc3 = "Doctors suggest that driving may cause increased stress and blood pressure."
```

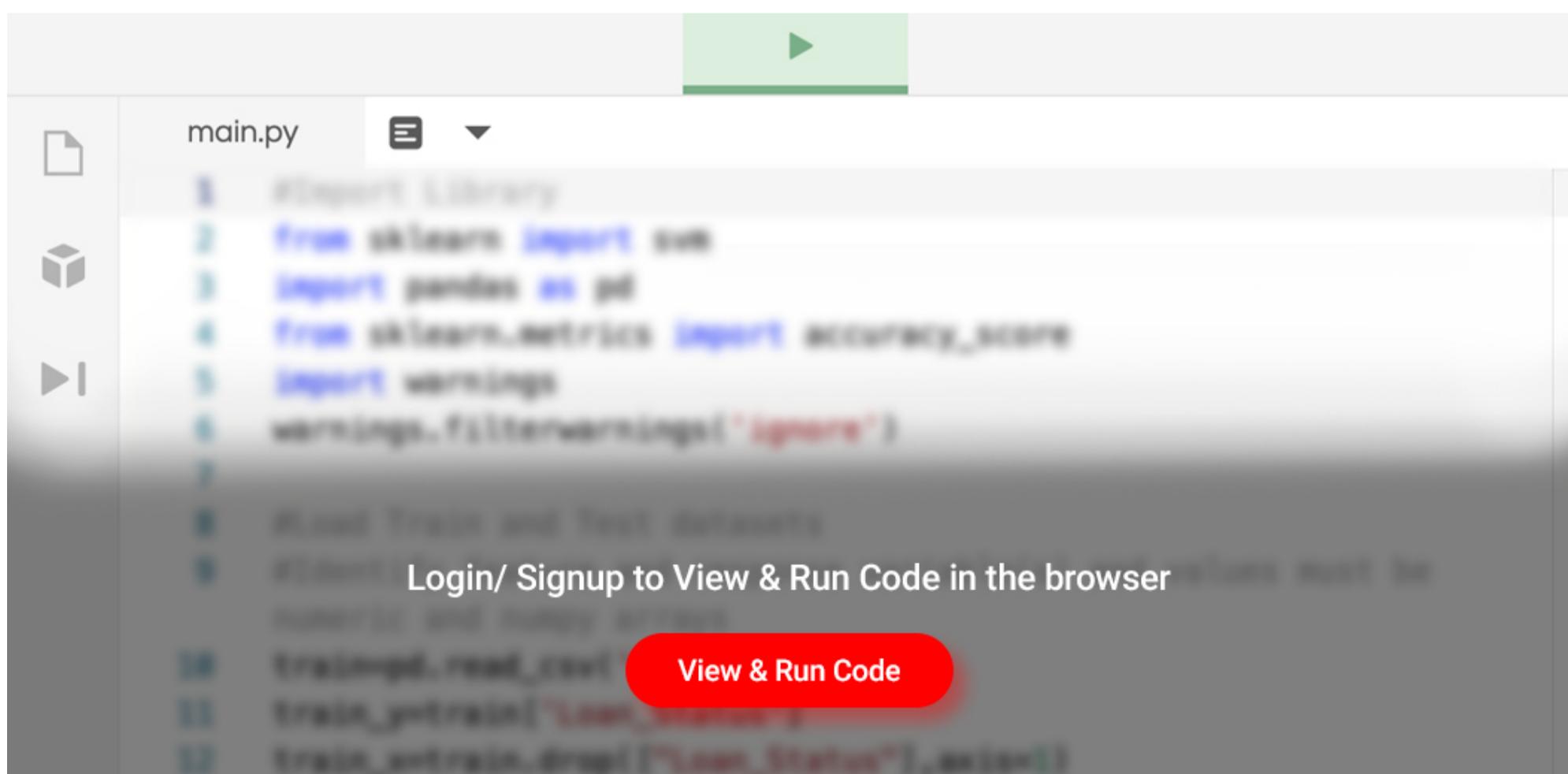
```
doc4 = "Sometimes I feel pressure to perform well at school, but my father never seems to drive my  
sister to do better."  
  
doc5 = "Health experts say that Sugar is not good for your lifestyle."  
  
# compile documents  
doc_complete = [doc1, doc2, doc3, doc4, doc5]
```

## Cleaning and Preprocessing

Cleaning is an important step before any text mining task, in this step, we will remove the punctuations, stopwords and normalize the corpus.

```

```
from nltk.corpus import stopwords  
  
from nltk.stem.wordnet import WordNetLemmatizer  
  
import string  
  
stop = set(stopwords.words('english'))  
  
exclude = set(string.punctuation)  
  
lemma = WordNetLemmatizer()  
  
def clean(doc):  
  
    stop_free = " ".join([i for i in doc.lower().split() if i not in stop])  
  
    punc_free = ''.join(ch for ch in stop_free if ch not in exclude)  
  
    normalized = " ".join(lemma.lemmatize(word) for word in punc_free.split())  
  
    return normalized  
  
doc_clean = [clean(doc).split() for doc in doc_complete]  
```
```



## Preparing Document-Term Matrix

All the text documents combined is known as the corpus. To run any mathematical model on text corpus, it is a good practice to convert it into a matrix representation. LDA model looks for repeating term patterns in the entire DT matrix. Python provides many great libraries for text mining practices, “gensim” is one such clean and beautiful library to handle text data. It is scalable, robust and efficient. Following code shows how to convert a corpus into a document-term matrix.

```
```
# Importing Gensim
import gensim
from gensim import corpora

# Creating the term dictionary of our courpus, where every unique term is assigned an index. dictionary
= corpora.Dictionary(doc_clean)

# Converting list of documents (corpus) into Document Term Matrix using dictionary prepared above.
doc_term_matrix = [dictionary.doc2bow(doc) for doc in doc_clean]
```
```

```

## Running LDA Model

Next step is to create an object for LDA model and train it on Document-Term matrix. The training also requires few parameters as input which are explained in the above section. The gensim module allows both LDA model estimation from a training corpus and inference of topic distribution on new, unseen documents.

```
```
# Creating the object for LDA model using gensim library
Lda = gensim.models.ldamodel.LdaModel

# Running and Trainign LDA model on the document term matrix.
ldamodel = Lda(doc_term_matrix, num_topics=3, id2word = dictionary, passes=50)
```
```

```

## Results

```
```
print(ldamodel.print_topics(num_topics=3, num_words=3))
['0.168*health + 0.083*sugar + 0.072*bad,
'0.061*consume + 0.050*drive + 0.050*sister,
'0.049*pressur + 0.049*father + 0.049*sister]
```
```

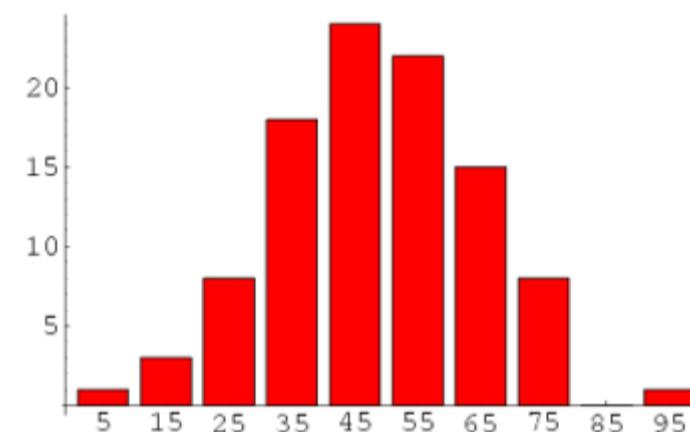
```

Each line is a topic with individual topic terms and weights. Topic1 can be termed as Bad Health, and Topic3 can be termed as Family.

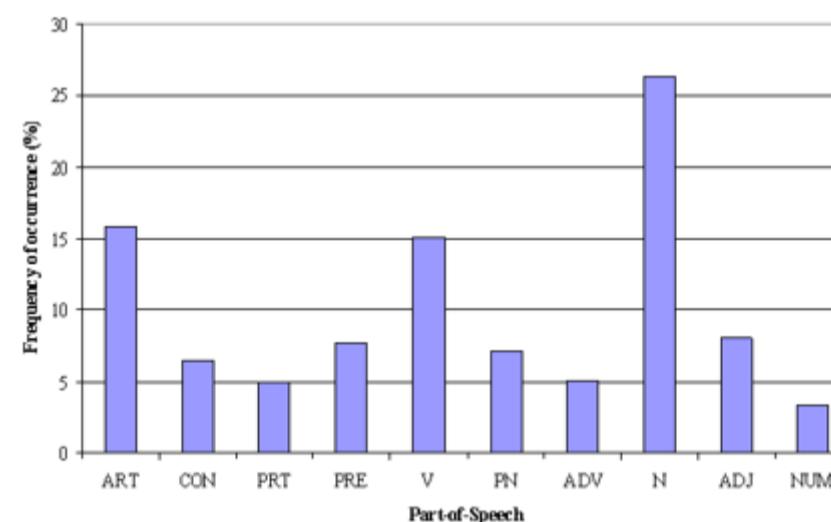
## Tips to improve results of topic modeling

The results of topic models are completely dependent on the features (terms) present in the corpus. The corpus is represented as document term matrix, which in general is very sparse in nature. Reducing the dimensionality of the matrix can improve the results of topic modelling. Based on my practical experience, there are few approaches which do the trick.

**1. Frequency Filter** – Arrange every term according to its frequency. Terms with higher frequencies are more likely to appear in the results as compared ones with low frequency. The low frequency terms are essentially weak features of the corpus, hence it is a good practice to get rid of all those weak features. An exploratory analysis of terms and their frequency can help to decide what frequency value should be considered as the threshold.



**2. Part of Speech Tag Filter** – POS tag filter is more about the context of the features than frequencies of features. Topic Modelling tries to map out the recurring patterns of terms into topics. However, every term might not be equally important contextually. For example, POS tag IN contain terms such as – “within”, “upon”, “except”. “CD” contains – “one”, “two”, “hundred” etc. “MD” contains “may”, “must” etc. These terms are the supporting words of a language and can be removed by studying their post tags.



**3. Batch Wise LDA** – In order to retrieve most important topic terms, a corpus can be divided into batches of fixed sizes. Running LDA multiple times on these batches will provide different results, however, the best topic terms will be the intersection of all batches.

Note: If you want to learn Topic Modeling in detail and also do a project using it, then we have a [video based course on NLP, covering Topic Modeling and its implementation in Python](#).

## Topic Modelling for Feature Selection

Sometimes LDA can also be used as feature selection technique. Take an example of text classification problem where the training data contain category wise documents. If LDA is running on sets of category wise documents. Followed by removing common topic terms across the results of different categories will give the best features for a category.

## Endnotes

With this, we come to this end of tutorial on Topic Modeling. I hope this will help you to improve your knowledge to work on text data. To reap maximum benefits out of this tutorial, I'd suggest you practice the codes side by side and check the results.

Did you find the article useful? Share with us if you have done similar kind of analysis before. Do let us know your thoughts about this article in the box below.

## References

- [http://link.springer.com/chapter/10.1007%2F978-3-642-13657-3\\_43](http://link.springer.com/chapter/10.1007%2F978-3-642-13657-3_43)

Got expertise in Business Intelligence / Machine Learning / Big Data / Data Science? Showcase your knowledge and help Analytics Vidhya community by [posting your blog](#).

---

[Advanced Python](#) [analytic modeling](#) [Analytics](#) [data science in python](#) [feature selection](#)

[latent dirichlet allocation](#) [live coding](#) [Natural language processing](#) [NLP](#) [topic modeling](#)

---

## About the Author



[Shivam5992 Bansal](#)

Shivam Bansal is a data scientist with exhaustive experience in Natural Language Processing and Machine Learning in several domains. He is passionate about learning and always looks forward to solving challenging analytical problems.

## Our Top Authors



[view more](#)

---

## Download

Analytics Vidhya App for the Latest blog/Article



[Previous Post](#)

[Bringing Analytics into Indian Film Industry with Back Tracing Algorithm](#)

[Next Post](#)

[10 Real World Applications of Internet of Things \(IoT\) – Explained in Videos](#)