

Machine Learning Plus
[\(<https://www.machinelearningplus.com/>\)](https://www.machinelearningplus.com/)

Topic modeling visualization – How to present the results of LDA models?

by Selva Prabhakaran (<https://www.machinelearningplus.com/author/selva86/>) | Posted on December 4, 2018 (<https://www.machinelearningplus.com/nlp/topic-modeling-visualization-how-to-present-results lda-models/>)

In this post, we discuss techniques to visualize the output and results from topic model (LDA) based on the gensim package.



Search ...

Search

Recent Posts

[data.table in R – The Complete Beginners Guide](https://www.machinelearningplus.com/data-manipulation/datatable-in-r-complete-guide/)
(<https://www.machinelearningplus.com/data-manipulation/datatable-in-r-complete-guide/>)

[Augmented Dickey Fuller Test \(ADF Test\) – Must Read Guide](https://www.machinelearningplus.com/time-series/augmented-dickey-fuller-test/)
(<https://www.machinelearningplus.com/time-series/augmented-dickey-fuller-test/>)

[KPSS Test for Stationarity](https://www.machinelearningplus.com/time-series/kpss-test-for-stationarity/)
(<https://www.machinelearningplus.com/time-series/kpss-test-for-stationarity/>)

[101 R data.table Exercises](https://www.machinelearningplus.com/data-manipulation/101-r-data-table-exercises/)
(<https://www.machinelearningplus.com/data-manipulation/101-r-data-table-exercises/>)

[P-Value – Understanding from Scratch](https://www.machinelearningplus.com/statistics/p-value/)
(<https://www.machinelearningplus.com/statistics/p-value/>)

[101 Python datatableViews \(pandasDataViews\)](#)



Topic modeling visualization – How to present the results of LDA models?

Contents

[columnize]

1. [Introduction](#)
2. [Import NewsGroups Dataset](#)
3. [Tokenize Sentences and Clean](#)
4. [Build the Bigram, Trigram Models and Lemmatize](#)
5. [Build the Topic Model](#)
- [Presenting the Results](#)

6. [What is the Dominant topic and its percentage contribution in each document?](#)
7. [The most representative sentences for each topic](#)
8. [Frequency Distribution of Word Counts in Documents](#)
9. [Word Clouds of Top N Keywords in Each Topic](#)
10. [Word Counts of Topic Keywords](#)
11. [Sentence Chart Colored by Topic](#)
12. [What are the most discussed topics in the documents?](#)
13. [t-SNE Clustering Chart](#)
14. [pyLDAVis](#)
15. [Conclusion](#)

[/columnize]

1. Introduction

In [topic modeling with gensim](#) (<https://www.machinelearningplus.com/nlp/topic-modeling-gensim-python/>), we followed a structured workflow to build an insightful topic model based on the Latent Dirichlet Allocation (LDA) algorithm.

In this post, we will build the topic model using gensim's native `LdaModel` and explore multiple strategies to effectively visualize the results using `matplotlib` plots.

I will be using a portion of the 20 Newsgroups dataset since the focus is more on approaches to visualizing the results.

Let's begin by importing the packages and the 20 News Groups dataset.

```
import sys
# !{sys.executable} -m spacy download en
import re, numpy as np, pandas as pd
from pprint import pprint

# Gensim
import gensim, spacy, logging, warnings
import gensim.corpora as corpora
from gensim.utils import lemmatize, simple_preprocess
from gensim.models import CoherenceModel
import matplotlib.pyplot as plt
```

[101 Python Datatable Exercises \(Pydatatable\)](#)
(<https://www.machinelearningplus.com/data-manipulation/101-python-datatable-exercises-pydatatable/>).

[Vector Autoregression \(VAR\) – Comprehensive Guide with Examples in Python](#)

(<https://www.machinelearningplus.com/time-series/vector-autoregression-examples-python/>).

[Mahalanobis Distance – Understanding the math with examples \(python\)](#)
(<https://www.machinelearningplus.com/statistics/mahalanobis-distance/>).

[datetime in Python – Simplified Guide with Clear Examples](#)
(<https://www.machinelearningplus.com/python/datetime-python-examples/>).

[Principal Component Analysis \(PCA\) – Better Explained](#)
(<https://www.machinelearningplus.com/machine-learning/principal-components-analysis-pca-better-explained/>).

[Python Logging – Simplest Guide with Full Code and Examples](#)
(<https://www.machinelearningplus.com/python/logging-guide/>).

Feedback

[Matplotlib Histogram – How to Visualize Distributions in Python](#)
(<https://www.machinelearningplus.com/plots/matplotlib-histogram-python-examples/>).

[ARIMA Model – Complete Guide to Time Series Forecasting in Python](#)
(<https://www.machinelearningplus.com/time-series/arima-model-time-series-forecasting-python/>).

[Time Series Analysis in Python – A Comprehensive Guide with Examples](#)
(<https://www.machinelearningplus.com/time-series/time-series-analysis-python/>).

[Matplotlib Tutorial – A Complete Guide to Python Plot w/ Examples](#)
(<https://www.machinelearningplus.com/plots/matplotlib-tutorial-complete-guide-python-plot-examples/>).

[Topic modeling visualization – How to present the results of LDA models?](#)
(<https://www.machinelearningplus.com/nlp/topic-modeling-visualization-how-to-present-results-lda-models/>).

[Top 50 matplotlib Visualizations – The Master Plots \(with full python code\)](#)
(<https://www.machinelearningplus.com/plots/top-50-matplotlib-visualizations-the-master-plots-python/>).

```

# NLTK Stop words
from nltk.corpus import stopwords
stop_words = stopwords.words('english')
stop_words.extend(['from', 'subject', 're', 'edu', 'use', 'not', 'would', 'say', 'could', '_',
                   '%matplotlib inline
warnings.filterwarnings("ignore",category=DeprecationWarning)
logging.basicConfig(format='%(asctime)s : %(levelname)s : %(message)s', level=logging.ERROR)

```

2. Import NewsGroups Dataset

Let's import the news groups dataset and retain only 4 of the target_names categories.

```

# Import Dataset
df = pd.read_json('https://raw.githubusercontent.com/selva86/datasets/master/news'
df = df.loc[df.target_names.isin(['soc.religion.christian', 'rec.sport.hockey', 'talk.politics.
print(df.shape) #> (2361, 3)
df.head()

```

		content	target	target_names
10	From: irwin@cmptrc.lonestar.org (Irwin Arnstei...		8	rec.motorcycles
10007	From: jet@netcom.Netcom.COM (J. Eric Townsend)...		8	rec.motorcycles
10008	From: gld@cunixb.cc.columbia.edu (Gary L Dare)...		10	rec.sport.hockey
10017	From: nstramer@supergas.dazixco.ingr.com (Naft...		17	talk.politics.mideast
10019	From: mussack@austin.ibm.com (Christopher Muss...		15	soc.religion.christian

(https://machinelearningplus.sirv.com/WP_www.machinelearningplus.com/2018/12/i/n/n/input_table-min.png).

3. Tokenize Sentences and Clean

Removing the emails, new line characters, single quotes and finally split the sentence into a list of words using gensim's simple_preprocess(). Setting the deacc=True option removes punctuations.

```

def sent_to_words(sentences):
    for sent in sentences:
        sent = re.sub('\S*@\S*\s?', '', sent) # remove emails
        sent = re.sub('\s+', ' ', sent) # remove newline chars
        sent = re.sub("\'", "", sent) # remove single quotes
        sent = gensim.utils.simple_preprocess(str(sent), deacc=True)
        yield(sent)

    # Convert to list
    data = df.content.values.tolist()
    data_words = list(sent_to_words(data))

```

[List Comprehensions in Python – My Simplified Guide](#)
[\(https://www.machinelearningplus.com/python/list-comprehensions-in-python/\)](https://www.machinelearningplus.com/python/list-comprehensions-in-python/).

[Python @Property Explained – How to Use and When? \(Full Examples\)](#)
[\(https://www.machinelearningplus.com/python/python-property/\)](https://www.machinelearningplus.com/python/python-property/).

[How Naive Bayes Algorithm Works? \(with example and full code\)](#)
[\(https://www.machinelearningplus.com/predictive-modeling/how-naive-bayes-algorithm-works-with-example-and-full-code/\)](https://www.machinelearningplus.com/predictive-modeling/how-naive-bayes-algorithm-works-with-example-and-full-code/).

Tags

[Bigrams](#) (<https://www.machinelearningplus.com/tag/bigrams/>)

[Classification](#)
[\(https://www.machinelearningplus.com/tag/classification/\)](https://www.machinelearningplus.com/tag/classification/)

[Corpus](#) (<https://www.machinelearningplus.com/tag/corpus/>)

[Cosine Similarity](#)
[\(https://www.machinelearningplus.com/tag/cosine-similarity/\)](https://www.machinelearningplus.com/tag/cosine-similarity/)

[data.table](#)
[\(https://www.machinelearningplus.com/tag/datatable/\)](https://www.machinelearningplus.com/tag/datatable/)

[Data Manipulation](#)
[\(https://www.machinelearningplus.com/tag/manipulation/\)](https://www.machinelearningplus.com/tag/manipulation/)

[Debugging](#)
[\(https://www.machinelearningplus.com/tag/debugging/\)](https://www.machinelearningplus.com/tag/debugging/)

[Doc2Vec](#)
[\(https://www.machinelearningplus.com/tag/doc2vec/\)](https://www.machinelearningplus.com/tag/doc2vec/)

[Feedback](#)

[Evaluation Metrics](#)
[\(https://www.machinelearningplus.com/tag/evaluation-metrics/\)](https://www.machinelearningplus.com/tag/evaluation-metrics/)

[FastText](#)
[\(https://www.machinelearningplus.com/tag/fasttext/\)](https://www.machinelearningplus.com/tag/fasttext/)

[Feature Selection](#)
[\(https://www.machinelearningplus.com/tag/feature-selection/\)](https://www.machinelearningplus.com/tag/feature-selection/)

[Gensim](#)
[\(https://www.machinelearningplus.com/tag/gensim/\)](https://www.machinelearningplus.com/tag/gensim/)

[IQA](#)

[LDA](#)
[\(https://www.machinelearningplus.com/tag/lda/\)](https://www.machinelearningplus.com/tag/lda/)

[Lemmatization](#)
[\(https://www.machinelearningplus.com/tag/lemmatization/\)](https://www.machinelearningplus.com/tag/lemmatization/)

[Linear Regression](#)
[\(https://www.machinelearningplus.com/tag/linear-regression/\)](https://www.machinelearningplus.com/tag/linear-regression/)

[Logistic](#)
[\(https://www.machinelearningplus.com/tag/logistic/\)](https://www.machinelearningplus.com/tag/logistic/)

[LSI](#)
[\(https://www.machinelearningplus.com/tag/lsi/\)](https://www.machinelearningplus.com/tag/lsi/)

[Matplotlib](#)
[\(https://www.machinelearningplus.com/tag/matplotlib/\)](https://www.machinelearningplus.com/tag/matplotlib/)

[Multiprocessing](#)
[\(https://www.machinelearningplus.com/tag/multiprocessing/\)](https://www.machinelearningplus.com/tag/multiprocessing/)

[NLP](#)
[\(https://www.machinelearningplus.com/tag/nlp/\)](https://www.machinelearningplus.com/tag/nlp/)

[NLTK](#)
[\(https://www.machinelearningplus.com/tag/nltk/\)](https://www.machinelearningplus.com/tag/nltk/)

[Numpy](#)
[\(https://www.machinelearningplus.com/tag\(numpy/\)](https://www.machinelearningplus.com/tag(numpy/)

[P-Value](#)
[\(https://www.machinelearningplus.com/tag/p-value/\)](https://www.machinelearningplus.com/tag/p-value/)

[Pandas](#)
[\(https://www.machinelearningplus.com/tag/pandas/\)](https://www.machinelearningplus.com/tag/pandas/)

[Parallel Processing](#)
[\(https://www.machinelearningplus.com/tag/parallel-processing/\)](https://www.machinelearningplus.com/tag/parallel-processing/)

[Phraser](#)

```

print(data_words[:1])
# [[['from', 'irwin', 'arnstein', 'subject', 're', 'recommendation', 'on', 'duc', 'summary', 'wh

```

[Practice](https://www.machinelearningplus.com/tag/phraser/)

[Exercise](https://www.machinelearningplus.com/tag/practice-exercise/)

Python

[\(https://www.machinelearningplus.com/\)](https://www.machinelearningplus.com/)

R

[\(https://www.machinelearningplus.com/\)](https://www.machinelearningplus.com/)

[Regex](https://www.machinelearningplus.com/tag/regex/)

[Regression](https://www.machinelearningplus.com/tag/regression/)

[Residual Analysis](https://www.machinelearningplus.com/tag/residual-analysis/)

[\(https://www.machinelearningplus.com/tag/residual-analysis/\)](https://www.machinelearningplus.com/tag/residual-analysis/)

[Scikit Learn](https://www.machinelearningplus.com/tag/scikit-learn/)

[\(https://www.machinelearningplus.com/tag/scikit-learn/\)](https://www.machinelearningplus.com/tag/scikit-learn/)

[Significance Tests](https://www.machinelearningplus.com/tag/significance-tests/)

[\(https://www.machinelearningplus.com/tag/significance-tests/\)](https://www.machinelearningplus.com/tag/significance-tests/)

[Soft Cosine Similarity](https://www.machinelearningplus.com/tag/soft-cosine-similarity/)

[\(https://www.machinelearningplus.com/tag/soft-cosine-similarity/\)](https://www.machinelearningplus.com/tag/soft-cosine-similarity/)

[spaCy](https://www.machinelearningplus.com/tag/spacy/)

[\(https://www.machinelearningplus.com/tag/spacy/\)](https://www.machinelearningplus.com/tag/spacy/)

[Stationarity](https://www.machinelearningplus.com/tag/stationarity/)

[\(https://www.machinelearningplus.com/tag/stationarity/\)](https://www.machinelearningplus.com/tag/stationarity/)

[Summarization](https://www.machinelearningplus.com/tag/summarization/)

[\(https://www.machinelearningplus.com/tag/summarization/\)](https://www.machinelearningplus.com/tag/summarization/)

[TaggedDocument](https://www.machinelearningplus.com/tag/taggeddocument/)

[\(https://www.machinelearningplus.com/tag/taggeddocument/\)](https://www.machinelearningplus.com/tag/taggeddocument/)

[TextBlob](https://www.machinelearningplus.com/tag/textblob/)

[\(https://www.machinelearningplus.com/tag/textblob/\)](https://www.machinelearningplus.com/tag/textblob/)

[TFIDF](https://www.machinelearningplus.com/tag/tfidf/)

[\(https://www.machinelearningplus.com/tag/tfidf/\)](https://www.machinelearningplus.com/tag/tfidf/)

[Time Series](https://www.machinelearningplus.com/tag/time-series/)

[\(https://www.machinelearningplus.com/tag/time-series/\)](https://www.machinelearningplus.com/tag/time-series/)

[Feedback](https://www.machinelearningplus.com/tag/feedback/)

[series/\)](https://www.machinelearningplus.com/tag/series/)

[Topic Modeling](https://www.machinelearningplus.com/tag/topic-modeling/)

[\(https://www.machinelearningplus.com/tag/topic-modeling/\)](https://www.machinelearningplus.com/tag/topic-modeling/)

[Visualization](https://www.machinelearningplus.com/tag/visualization/)

[\(https://www.machinelearningplus.com/tag/visualization/\)](https://www.machinelearningplus.com/tag/visualization/)

[Word2Vec](https://www.machinelearningplus.com/tag/word2vec/)

[\(https://www.machinelearningplus.com/tag/word2vec/\)](https://www.machinelearningplus.com/tag/word2vec/)

```

# Build the bigram and trigram models
bigram = gensim.models.Phrases(data_words, min_count=5, threshold=100) # higher threshold fewer
trigram = gensim.models.Phrases(bigram[data_words], threshold=100)
bigram_mod = gensim.models.phrases.Phraser(bigram)
trigram_mod = gensim.models.phrases.Phraser(trigram)

# !python3 -m spacy download en # run in terminal once
# or do
# !conda install -c conda-forge spacy-model-en_core_web_md
# and use nlp=spacy.load('en_core_web_sm') instead in below function.

def process_words(texts, stop_words=stop_words, allowed_postags=['NOUN', 'ADJ', 'VERB', 'ADV']):
    """Remove Stopwords, Form Bigrams, Trigrams and Lemmatization"""
    texts = [[word for word in simple_preprocess(str(doc)) if word not in stop_words] for doc in i
    texts = [bigram_mod[doc] for doc in texts]
    texts = [trigram_mod[bigram_mod[doc]] for doc in texts]
    texts_out = []
    nlp = spacy.load('en', disable=['parser', 'ner'])
    for sent in texts:
        doc = nlp(" ".join(sent))
        texts_out.append([token.lemma_ for token in doc if token.pos_ in allowed_postags])
    # remove stopwords once more after lemmatization
    texts_out = [[word for word in simple_preprocess(str(doc)) if word not in stop_words] for d
    return texts_out

data_ready = process_words(data_words) # processed Text Data!

```

4. Build the Bigram, Trigram Models and Lemmatize

Let's form the bigram and trigrams using the Phrases model. This is passed to Phraser() for efficiency in speed of execution.

Next, lemmatize each word to its root form, keeping only nouns, adjectives, verbs and adverbs.

We keep only these POS tags because they are the ones contributing the most to the meaning of the sentences. Here, I use spacy for [lemmatization](https://www.machinelearningplus.com/nlp/lemmatization-examples-python/) (<https://www.machinelearningplus.com/nlp/lemmatization-examples-python/>).

5. Build the Topic Model

To build the LDA topic model using `LdaModel()`, you need the corpus and the dictionary. Let's create them first and then build the model. The trained topics (keywords and weights) are printed below as well.

If you examine the topic key words, they are nicely segregate and collectively represent the topics we initially chose: Christianity, Hockey, MidEast and Motorcycles. Nice!

[Feedback](#)

```
# Create Dictionary
id2word = corpora.Dictionary(data_ready)

# Create Corpus: Term Document Frequency
corpus = [id2word.doc2bow(text) for text in data_ready]

# Build LDA model
lda_model = gensim.models.ldamodel.LdaModel(corpus=corpus,
                                              id2word=id2word,
                                              num_topics=4,
                                              random_state=100,
                                              update_every=1,
                                              chunksize=10,
                                              passes=10,
                                              alpha='symmetric',
                                              iterations=100,
                                              per_word_topics=True)

pprint(lda_model.print_topics())
#> [(0,
#>   '0.017*"write" + 0.015*"people" + 0.014*"organization" + 0.014*"article" + '
#>   '0.013*"time" + 0.008*"give" + 0.008*"first" + 0.007*"tell" + 0.007*"new" + '
#>   '0.007*"question"' ),
#>   (1,
#>   '0.008*"christian" + 0.008*"believe" + 0.007*"god" + 0.007*"Law" + '
#>   '0.006*"state" + 0.006*"israel" + 0.006*"israeli" + 0.005*"exist" + '
#>   '0.005*"way" + 0.004*"bible"' ),
#>   (2,
```

```

#> '0.024*"armenian" + 0.012*"bike" + 0.006*"kill" + 0.006*"work" +
#> '0.005*"well" + 0.005*"year" + 0.005*"sumgait" + 0.005*"soldier" +
#> '0.004*"way" + 0.004*"ride"),
#> (3,
#> '0.019*"team" + 0.019*"game" + 0.013*"hockey" + 0.010*"player" +
#> '0.009*"play" + 0.009*"win" + 0.009*"nhl" + 0.009*"year" + 0.009*"hawk" +
#> '0.009*"season")]

```

6. What is the Dominant topic and its percentage contribution in each document

In LDA models, each document is composed of multiple topics. But, typically only one of the topics is dominant. The below code extracts this dominant topic for each sentence and shows the weight of the topic and the keywords in a nicely formatted output.

This way, you will know which document belongs predominantly to which topic.

[Feedback](#)

```

def format_topics_sentences(ldamodel=None, corpus=corpus, texts=data):
    # Init output
    sent_topics_df = pd.DataFrame()

    # Get main topic in each document
    for i, row_list in enumerate(ldamodel[corpus]):
        row = row_list[0] if ldamodel.per_word_topics else row_list
        # print(row)
        row = sorted(row, key=lambda x: (x[1]), reverse=True)
        # Get the Dominant topic, Perc Contribution and Keywords for each document
        for j, (topic_num, prop_topic) in enumerate(row):
            if j == 0: # => dominant topic
                wp = ldamodel.show_topic(topic_num)
                topic_keywords = ", ".join([word for word, prop in wp])
                sent_topics_df = sent_topics_df.append(pd.Series([int(topic_num), round(prop_topic, 4)], index=['Dominant_Topic', 'Perc_Contribution', 'Topic_Keywords']))
            else:
                break
    sent_topics_df.columns = ['Dominant_Topic', 'Perc_Contribution', 'Topic_Keywords']

    # Add original text to the end of the output
    contents = pd.Series(texts)
    sent_topics_df = pd.concat([sent_topics_df, contents], axis=1)
    return(sent_topics_df)

df_topic_sents_keywords = format_topics_sentences(ldamodel=lda_model, corpus=corpus, texts=data
# Format

```

```

df_dominant_topic = df_topic_sents_keywords.reset_index()
df_dominant_topic.columns = ['Document_No', 'Dominant_Topic', 'Topic_Perc_Contrib', 'Keywords',
df_dominant_topic.head(10)

```

Document_No	Dominant_Topic	Topic_Perc_Contrib	Keywords	Text
0	0	2.0	0.6916 armenian, bike, kill, work, well, year, sumgai...	[irwin, arnstein, recommendation, duc, summary...
1	1	2.0	0.6247 armenian, bike, kill, work, well, year, sumgai...	[eric, townsend, insurance, lotsa_point, reply...
2	2	0.0	0.5619 write, people, organization, article, time, gl...	[gary_dare, abc_coverage, reply, gary_dare, or...
3	3	1.0	0.4829 christian, believe, god, law, state, israel, i...	[naftaly_stramer, peace_talk, reply, organizat...
4	4	0.0	0.8043 write, people, organization, article, time, gl...	[question, authority, trendy, liberal, feminis...
5	5	1.0	0.5241 christian, believe, god, law, state, israel, i...	[harrass, work, prayer_reply, organization, u...
6	6	0.0	0.6817 write, people, organization, article, time, gl...	[andrew_infante, ok, little_hasty, organizatio...
7	7	1.0	0.6250 christian, believe, god, law, state, israel, i...	[james_sledd, afterlife, organization, social...
8	8	3.0	0.5280 team, game, hockey, player, play, win, nhl, ye...	[canadian, stanley_cup, organization, canada, ...
9	9	3.0	0.8933 team, game, hockey, player, play, win, nhl, ye...	[steve_gallicchio, possible, canadian, wc, team...

(https://machinelearningplus.sirv.com/WP www.machinelearningplus.com/2018/12/d/o/n/dominant_topics-min.png).

7. The most representative sentence for each topic

Sometimes you want to get samples of sentences that most represent a given topic. This code gets the most exemplar sentence for each topic.

[Feedback](#)

```

# Display setting to show more characters in column
pd.options.display.max_colwidth = 100

sent_topics_sorteddf_mallet = pd.DataFrame()
sent_topics_outdf_grpd = df_topic_sents_keywords.groupby('Dominant_Topic')

for i, grp in sent_topics_outdf_grpd:
    sent_topics_sorteddf_mallet = pd.concat([sent_topics_sorteddf_mallet,
                                              grp.sort_values(['Perc_Contribution'], ascending=False).head(1),
                                              axis=0])

# Reset Index
sent_topics_sorteddf_mallet.reset_index(drop=True, inplace=True)

# Format
sent_topics_sorteddf_mallet.columns = ['Topic_Num', "Topic_Perc_Contrib", "Keywords", "Representative_Text"]

# Show
sent_topics_sorteddf_mallet.head(10)

```

Topic_Num	Topic_Perc_Contrib	Keywords	Representative Text
0	0.0	0.9740 write, people, organization, article, time, give, first, tell, new, question	[zionism, racism, vn_mcgill_ca, organization, mcgill, university, author, write, thought, zionis...
1	1.0	0.8662 christian, believe, god, law, state, israel, israeli, exist, way, bible	[itm_clock, final_solution, gaza, oac, uci, organization, university, article, center, policy, res...
2	2.0	0.9600 armenian, bike, kill, work, well, year, sumgai, soldier, way, ride	[jhris, pith, live, free, quietly, die, organization, microsoft, corporation, rec_motorcycle, w...
3	3.0	0.9984 team, game, hockey, player, play, win, nhl, year, hawk, season	[result, game, play, sit, april, cook_charlie, organization, university, new_brunswick, tampa_ba...

(https://machinelearningplus.sirv.com/WP www.machinelearningplus.com/2018/12/m/o/n/most_representative_sentence-min.png).

8. Frequency Distribution of Word Counts in Documents

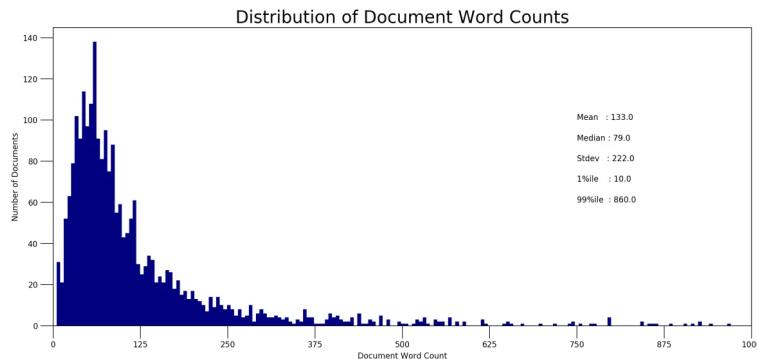
When working with a large number of documents, you want to know how big the documents are as a whole and by topic. Let's plot the document word counts distribution.

```
doc_lens = [len(d) for d in df_dominant_topic.Text]

# Plot
plt.figure(figsize=(16,7), dpi=160)
plt.hist(doc_lens, bins = 1000, color='navy')
plt.text(750, 100, "Mean : " + str(round(np.mean(doc_lens))))
plt.text(750, 90, "Median : " + str(round(np.median(doc_lens))))
plt.text(750, 80, "Stdev : " + str(round(np.std(doc_lens))))
plt.text(750, 70, "1%ile : " + str(round(np.quantile(doc_lens, q=0.01))))
plt.text(750, 60, "99%ile : " + str(round(np.quantile(doc_lens, q=0.99)))

plt.gca().set(xlim=(0, 1000), ylabel='Number of Documents', xlabel='Document Word Count')
plt.tick_params(size=16)
plt.xticks(np.linspace(0,1000,9))
plt.title('Distribution of Document Word Counts', fontdict=dict(size=22))
plt.show()
```

[Feedback](#)



(https://machinelearningplus.sirv.com/WP www.machinelearningplus.com/2018/12/1/n/01_distribution_of_document_wordcounts-min.png).

```
import seaborn as sns
import matplotlib.colors as mcolors
cols = [color for name, color in mcolors.TABLEAU_COLORS.items()] # more colors: 'mcolors.XKCD_'

fig, axes = plt.subplots(2,2, figsize=(16,14), dpi=160, sharex=True, sharey=True)

for i, ax in enumerate(axes.flatten()):
    df_dominant_topic_sub = df_dominant_topic.loc[df_dominant_topic.Dominant_Topic == i, :]
    doc_lens = [len(d) for d in df_dominant_topic_sub.Text]
    ax.hist(doc_lens, bins = 1000, color=cols[i])
    ax.tick_params(axis='y', labelcolor=cols[i], color=cols[i])
    sns.kdeplot(doc_lens, color="black", shade=False, ax=ax.twinx())
    ax.set(xlim=(0, 1000), xlabel='Document Word Count')
    ax.set_ylabel('Number of Documents', color=cols[i])
```

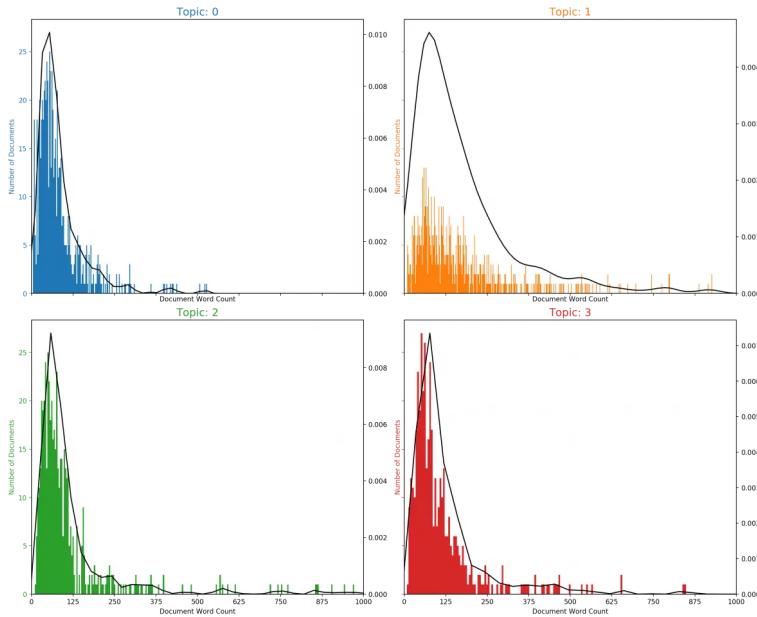
```

ax.set_title('Topic: ' + str(i), fontweight='bold', size=16, color=cols[i])

fig.tight_layout()
fig.subplots_adjust(top=0.90)
plt.xticks(np.linspace(0,1000,9))
fig.suptitle('Distribution of Document Word Counts by Dominant Topic', fontsize=22)
plt.show()

```

Distribution of Document Word Counts by Dominant Topic



Feedback

(https://machinelearningplus.sirv.com/WP_www.machinelearningplus.com/2018/12/2/n/02_distribution_of_document_wordcounts_by_topic_min.png).

9. Word Clouds of Top N Keywords in Each Topic

Though you've already seen what are the topic keywords in each topic, a word cloud with the size of the words proportional to the weight is a pleasant sight. The coloring of the topics I've taken here is followed in the subsequent plots as well.

```

# 1. WordCloud of Top N words in each topic
from matplotlib import pyplot as plt
from wordcloud import WordCloud, STOPWORDS
import matplotlib.colors as mcolors

cols = [color for name, color in mcolors.TABLEAU_COLORS.items()] # more colors: 'mcolors.XKCD_
# colors

cloud = WordCloud(stopwords=stop_words,
                  background_color='white',
                  width=2500,
                  height=1800,
                  max_words=10,
                  colormap='tab10',
                  color_func=lambda *args, **kwargs: cols[i],
                  prefer_horizontal=1.0)

topics = lda_model.show_topics(formatted=False)

```

```

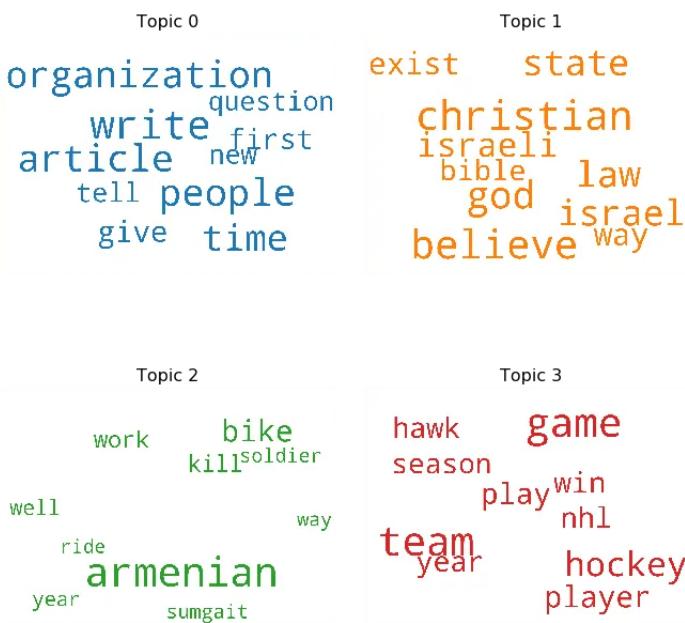
fig, axes = plt.subplots(2, 2, figsize=(10,10), sharex=True, sharey=True)

for i, ax in enumerate(axes.flatten()):
    fig.add_subplot(ax)
    topic_words = dict(topics[i][1])
    cloud.generate_from_frequencies(topic_words, max_font_size=300)
    plt.gca().imshow(cloud)
    plt.gca().set_title('Topic ' + str(i), fontdict=dict(size=16))
    plt.gca().axis('off')

plt.subplots_adjust(wspace=0, hspace=0)
plt.axis('off')
plt.margins(x=0, y=0)
plt.tight_layout()
plt.show()

```

[Feedback](#)



(https://machinelearningplus.sirv.com/WP_www.machinelearningplus.com/2018/12/3/n/03_topic_wordcloud-min.png)

10. Word Counts of Topic Keywords

When it comes to the keywords in the topics, the importance (weights) of the keywords matters. Along with that, how frequently the words have appeared in the documents is also interesting to look.

Let's plot the word counts and the weights of each keyword in the same chart.

You want to keep an eye out on the words that occur in multiple topics and the ones whose relative frequency is more than the weight. Often such words turn out to be less important. The chart I've drawn below is a result of adding several such words to the stop words list in the beginning and re-running the training process.

[Feedback](#)

```
from collections import Counter
topics = lda_model.show_topics(formatted=False)
data_flat = [w for w_list in data_ready for w in w_list]
counter = Counter(data_flat)

out = []
for i, topic in topics:
    for word, weight in topic:
        out.append([word, i, weight, counter[word]])

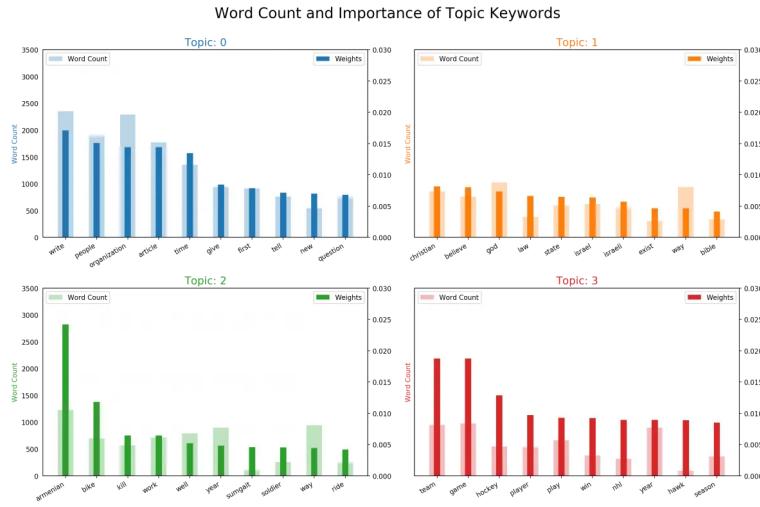
df = pd.DataFrame(out, columns=['word', 'topic_id', 'importance', 'word_count'])

# Plot Word Count and Weights of Topic Keywords
fig, axes = plt.subplots(2, 2, figsize=(16,10), sharey=True, dpi=160)
cols = [color for name, color in mcolors.TABLEAU_COLORS.items()]

for i, ax in enumerate(axes.flatten()):
    ax.bar(x='word', height="word_count", data=df.loc[df.topic_id==i, :], color=cols[i], width=1)
    ax_twin = ax.twinx()
    ax_twin.bar(x='word', height="importance", data=df.loc[df.topic_id==i, :], color=cols[i], width=1)
    ax.set_ylabel('Word Count', color=cols[i])
    ax_twin.set_ylabel('Importance', color=cols[i])
    ax.set_title('Topic: ' + str(i), color=cols[i], fontsize=16)
    ax.tick_params(axis='y', left=False)
    ax.set_xticklabels(df.loc[df.topic_id==i, 'word'], rotation=30, horizontalalignment='right')
    ax.legend(loc='upper left'); ax_twin.legend(loc='upper right')

fig.tight_layout(w_pad=2)
fig.suptitle('Word Count and Importance of Topic Keywords', fontsize=22, y=1.05)
```

```
plt.show()
```



(<https://machinelearningplus.sirv.com/WP www.machinelearningplus.com/2018/12/4/n/04 word count and importance of topics-min.png>).

11. Sentence Chart Colored by Topic

Each word in the document is representative of one of the 4 topics. Let's color each word in the given documents by the topic id it is attributed to.

The color of the enclosing rectangle is the topic assigned to the document.

[Feedback](#)

```
# Sentence Coloring of N Sentences
from matplotlib.patches import Rectangle

def sentences_chart(lda_model=lda_model, corpus=corpus, start = 0, end = 13):
    corp = corpus[start:end]
    mycolors = [color for name, color in mcolors.TABLEAU_COLORS.items()]

    fig, axes = plt.subplots(end-start, 1, figsize=(20, (end-start)*0.95), dpi=160)
    axes[0].axis('off')

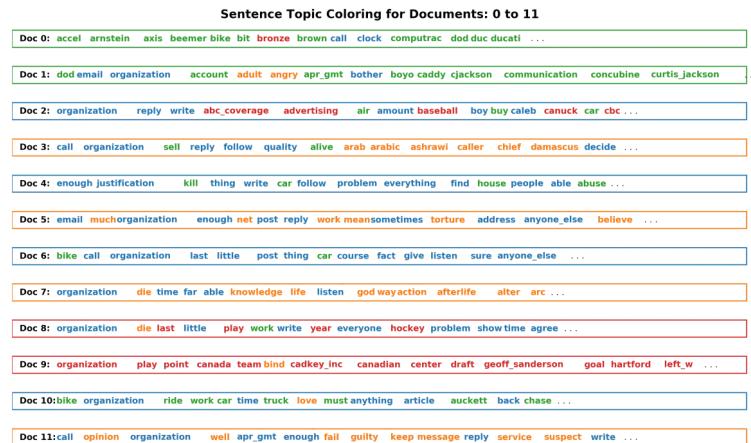
    for i, ax in enumerate(axes):
        if i > 0:
            corp_cur = corp[i-1]
            topic_percs, wordid_phivalue = lda_model[corp_cur]
            word_dominanttopic = [(lda_model.id2word[wd], topic[0]) for wd, topic in wordid_phivalue]
            ax.text(0.01, 0.5, "Doc " + str(i-1) + ": ", verticalalignment='center',
                    fontsize=16, color='black', transform=ax.transAxes, fontweight=700)

        # Draw Rectangle
        topic_percs_sorted = sorted(topic_percs, key=lambda x: (x[1]), reverse=True)
        ax.add_patch(Rectangle((0.0, 0.05), 0.99, 0.90, fill=None, alpha=1,
                              color=mycolors[topic_percs_sorted[0][0]], linewidth=2))

        word_pos = 0.06
        for j, (word, topics) in enumerate(word_dominanttopic):
            if j < 14:
                ax.text(word_pos, 0.5, word,
                        horizontalalignment='left',
                        verticalalignment='center',
```

```
    fontsize=16, color=mycolors[topics],  
    transform=ax.transAxes, fontweight=700)  
  
    word_pos += .009 * len(word) # to move the word for the next iter  
    ax.axis('off')  
  
    ax.text(word_pos, 0.5, '... ',  
            horizontalalignment='left',  
            verticalalignment='center',  
            fontsize=16, color='black',  
            transform=ax.transAxes)  
  
plt.subplots_adjust(wspace=0, hspace=0)  
plt.suptitle('Sentence Topic Coloring for Documents: ' + str(start) + ' to ' + str(end-2),  
plt.tight_layout()  
plt.show()  
  
sentences.chart()
```

Feedback



(https://machinelearningplus.sirv.com/WP_www.machinelearningplus.com/2018/12/5/n/05_sentence_coloring_chart-min.png).

12. What are the most discussed topics in the documents?

Let's compute the total number of documents attributed to each topic.

```
# Sentence Coloring of N Sentences

def topics_per_document(model, corpus, start=0, end=1):

    corpus_sel = corpus[start:end]
    dominant_topics = []
    topic_percentages = []
```

```

for i, corp in enumerate(corpus_sel):
    topic_percts, wordid_topics, wordid_phivalue = model[corp]
    dominant_topic = sorted(topic_percts, key = lambda x: x[1], reverse=True)[0][0]
    dominant_topics.append((i, dominant_topic))
    topic_percentages.append(topic_percts)

return(dominant_topics, topic_percentages)

dominant_topics, topic_percentages = topics_per_document(model=lda_model, corpus=corpus, end=-1

# Distribution of Dominant Topics in Each Document
df = pd.DataFrame(dominant_topics, columns=['Document_Id', 'Dominant_Topic'])
dominant_topic_in_each_doc = df.groupby('Dominant_Topic').size()
df_dominant_topic_in_each_doc = dominant_topic_in_each_doc.to_frame(name='count').reset_index()

# Total Topic Distribution by actual weight
topic_weightage_by_doc = pd.DataFrame([dict(t) for t in topic_percentages])
df_topic_weightage_by_doc = topic_weightage_by_doc.sum().to_frame(name='count').reset_index()

# Top 3 Keywords for each Topic
topic_top3words = [(i, topic) for i, topics in lda_model.show_topics(formatted=False)
                   for j, (topic, wt) in enumerate(topics) if j < 3]

df_top3words_stacked = pd.DataFrame(topic_top3words, columns=['topic_id', 'words'])
df_top3words = df_top3words_stacked.groupby('topic_id').agg(', \n'.join)
df_top3words.reset_index(inplace=True)

```

Let's make two plots:

[Feedback](#)

1. The number of documents for each topic by assigning the document to the topic that has the most weight in that document.
2. The number of documents for each topic by summing up the actual weight contribution of each topic to respective documents.

```

from matplotlib.ticker import FuncFormatter

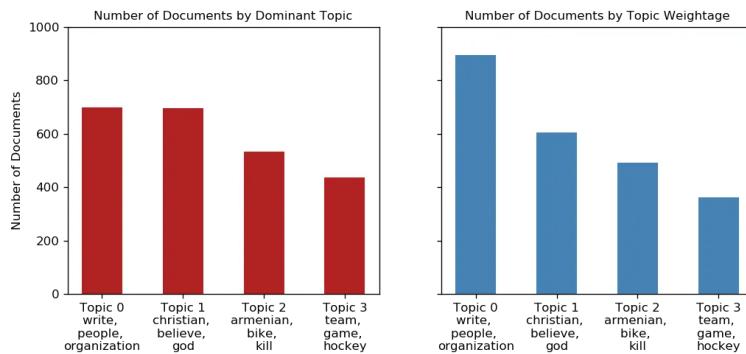
# Plot
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(10, 4), dpi=120, sharey=True)

# Topic Distribution by Dominant Topics
ax1.bar(x='Dominant_Topic', height='count', data=df_dominant_topic_in_each_doc, width=.5, color='red')
ax1.set_xticks(range(df_dominant_topic_in_each_doc.Dominant_Topic.unique().__len__()))
tick_formatter = FuncFormatter(lambda x, pos: 'Topic ' + str(x) + '\n' + df_top3words.loc[df_top3words['topic_id'] == x].words.iat[0])
ax1.xaxis.set_major_formatter(tick_formatter)
ax1.set_title('Number of Documents by Dominant Topic', fontdict=dict(size=10))
ax1.set_ylabel('Number of Documents')
ax1.set_yscale(0, 1000)

# Topic Distribution by Topic Weights
ax2.bar(x='index', height='count', data=df_topic_weightage_by_doc, width=.5, color='steelblue')
ax2.set_xticks(range(df_topic_weightage_by_doc.index.unique().__len__()))
ax2.xaxis.set_major_formatter(tick_formatter)
ax2.set_title('Number of Documents by Topic Weightage', fontdict=dict(size=10))

plt.show()

```



(https://machinelearningplus.sirv.com/WP www.machinelearningplus.com/2018/12/6/n/06_distribution_of_topics-min.png).

13. t-SNE Clustering Chart

Let's visualize the clusters of documents in a 2D space using t-SNE (t-distributed stochastic neighbor embedding) algorithm.

[Feedback](#)

```
# Get topic weights and dominant topics -----
from sklearn.manifold import TSNE
from bokeh.plotting import figure, output_file, show
from bokeh.models import Label
from bokeh.io import output_notebook

# Get topic weights
topic_weights = []
for i, row_list in enumerate(lda_model[corpus]):
    topic_weights.append([w for i, w in row_list[0]])

# Array of topic weights
arr = pd.DataFrame(topic_weights).fillna(0).values

# Keep the well separated points (optional)
arr = arr[np.amax(arr, axis=1) > 0.35]

# Dominant topic number in each doc
topic_num = np.argmax(arr, axis=1)

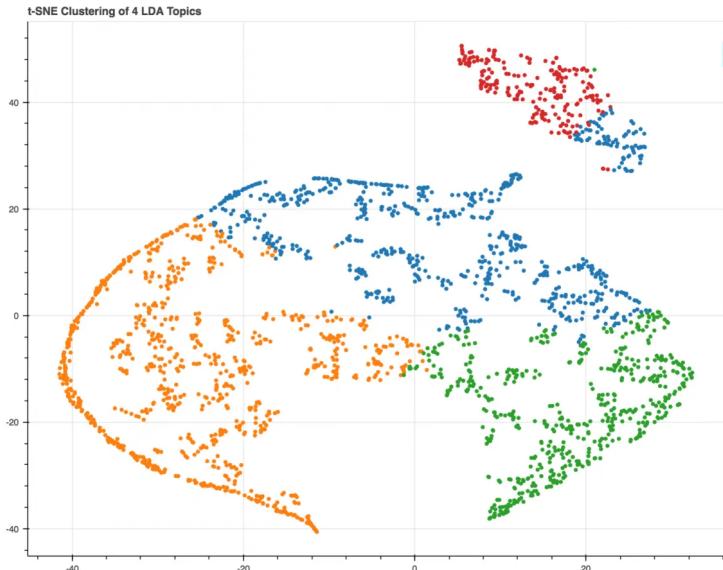
# tSNE Dimension Reduction
tsne_model = TSNE(n_components=2, verbose=1, random_state=0, angle=.99, init='pca')
tsne_lda = tsne_model.fit_transform(arr)

# Plot the Topic Clusters using Bokeh
output_notebook()
n_topics = 4
mycolors = np.array([color for name, color in mcolors.TABLEAU_COLORS.items()])
```

```

plot = figure(title="t-SNE Clustering of {} LDA Topics".format(n_topics),
              plot_width=900, plot_height=700)
plot.scatter(x=tsne_lda[:,0], y=tsne_lda[:,1], color=mycolors[topic_num])
show(plot)

```



(https://machinelearningplus.sirv.com/WP_www.machinelearningplus.com/2018/12/t/S/n/tSNE_Clustering-min.png).

14. pyLDAVis

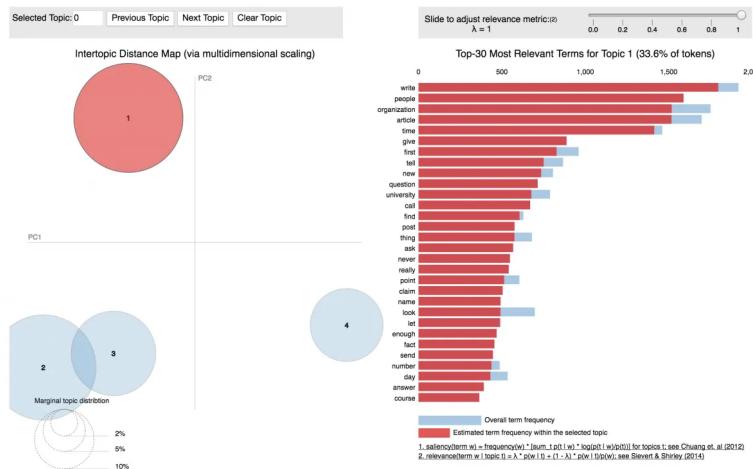
[Feedback](#)

Finally, pyLDAVis is the most commonly used and a nice way to visualise the information contained in a topic model. Below is the implementation for `LdaModel()`.

```

import pyLDAvis.gensim
pyLDAvis.enable_notebook()
vis = pyLDAvis.gensim.prepare(lda_model, corpus, dictionary=lda_model.id2word)
vis

```



(https://machinelearningplus.sirv.com/WP_www.machinelearningplus.com/2018/12/p/y/n/pyLDAVis-min.png).

15. Conclusion

13. CONCLUSION

We started from scratch by importing, cleaning and processing the newsgroups dataset to build the LDA model. Then we saw multiple ways to visualize the outputs of topic models including the word clouds and sentence coloring, which intuitively tells you what topic is dominant in each topic. A t-SNE clustering and the pyLDAVis are provide more details into the clustering of the topics.

Where next? If you are familiar with scikit learn, you can [build and grid search topic models using scikit learn](#) (<https://www.machinelearningplus.com/nlp/topic-modeling-python-sklearn-examples/>) as well.

What do you think?
30 Responses

 Upvote  Love

Comments Community  Privacy Policy  Login

 Recommend 2 Sort by Newest

Join the discussion...

LOG IN WITH

OR SIGN UP WITH DISQUS 

Name Feedback

 **andrea vega** • a month ago
Hi, I really hope that someone can help me.
I followed the example, with a little ups and downs but it worked until the last part,

see, in pyLDAVis, returns a text result but not the grafics, is there something that I missed?
p.d... working on Spyder (from anaconda)

thanks
^ | v • Reply •

 **Giuseppe** • a month ago
i can change the number of topics for the topic modelling, but how i change the number of topics for all the other analysis? In the Distribution of Document Word counts by Dominant topic, i always have 4 topics, how i can change the number?
Thanks
^ | v • Reply •

 **Abdullah Farooqi** • 2 months ago
Thanks a lot sir, your article helped me a lot in my final year project. But sir i have a problem in (pyLDAVis) portion, this says "no module of pyLDAVis named or found". please help me in this case
^ | v • Reply •

 **Felicia Fryer** • 3 months ago
Great article! I was able to reproduce your code which it has been very helpful in my journey of learning NLP.
Are you on LinkedIn? I will like to connect with you.
^ | v • Reply •