

A Beginner's Guide to Latent Dirichlet Allocation(LDA)

A statistical model for discovering the abstract topics aka topic modeling.



Ria Kulshrestha [Follow](#)

Jul 20, 2019 · 8 min read





Photo by Giulia Bertelli on Unsplash

What is topic modeling?

Topic modeling is a method for *unsupervised* classification of documents, similar to clustering on numeric data, which finds some natural groups of items (topics) even when we're not sure what we're looking for.

A document can be a part of multiple topics, kind of like in fuzzy clustering(soft clustering) in which each data point belongs to more than one cluster.

Why topic modeling?

Topic modeling provides methods for automatically organizing, understanding, searching, and summarizing large electronic archives. It can help with the following:

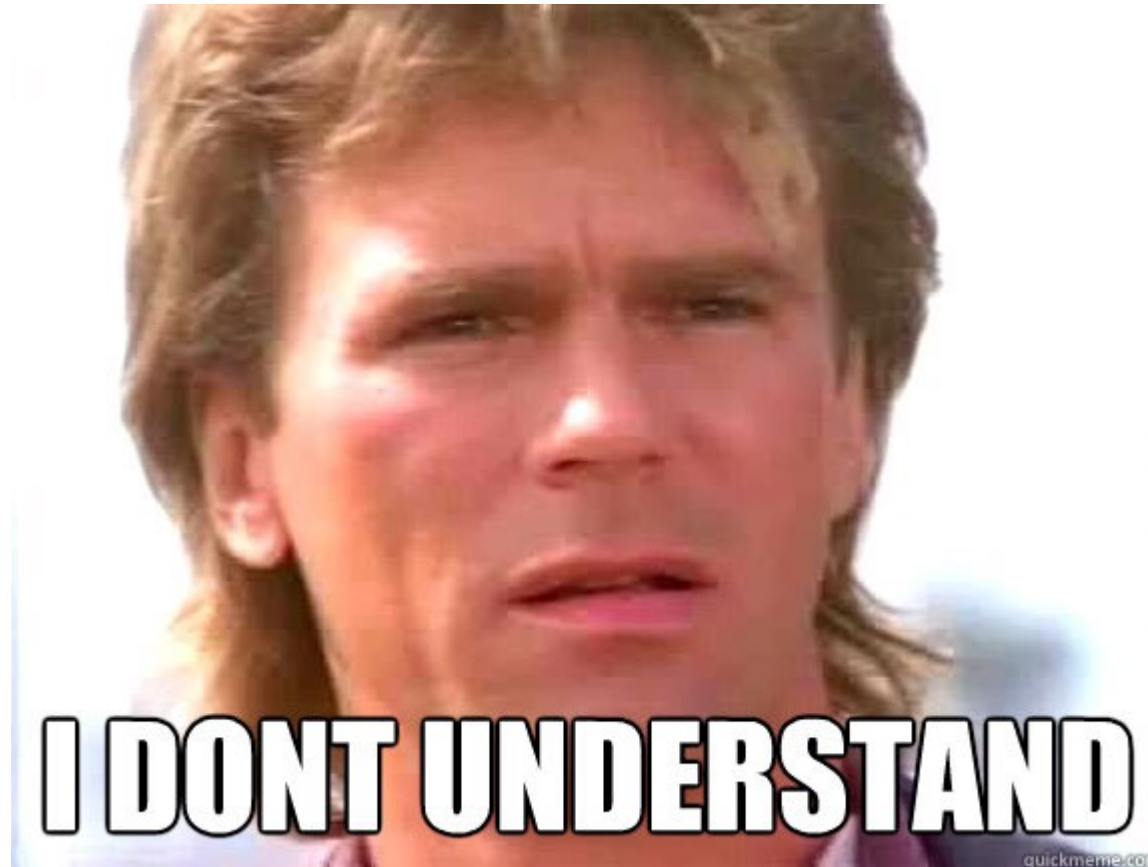
- discovering the hidden themes in the collection.
- classifying the documents into the discovered themes.
- using the classification to organize/summarize/search the documents.

For example, let's say a document belongs to the topics *food*, *dogs* and *health*. So if a user queries “*dog food*”, they might find the above-mentioned document relevant because it covers those topics(among other topics). We are able to figure its relevance with respect to the query without even going through the entire document.

Therefore, by annotating the document, based on the topics predicted by the modeling method, we are able to optimize our search process.

LDA

It is one of the most popular topic modeling methods. Each document is made up of various words, and each topic also has various words belonging to it. The aim of LDA is to find topics a document belongs to, based on the words in it. Confused much? Here is an example to walk you through it.



Model definition

Doc1: word1, word3, word5, word45, word11, word 62, word88 ...
 Doc2: word9, word77, word31, word58, word83, word 92, word49 ...
 Doc3: word44, word18, word52, word36, word64, word 11, word20 ...
 Doc4: word85, word62, word19, word4, word30, word 94, word67 ...
 Doc5: word19, word53, word74, word79, word45, word 39, word54 ...

Each document is a collection of words.

We have 5 documents each containing the words listed in front of them(ordered by frequency of occurrence).

What we want to figure out are the words in different topics, as shown in the table below. Each row in the table represents a different topic and each column a different word in the corpus. Each cell contains the probability that the word(column) belongs to the topic(row).

	Word1	word2	word3	word4
Topic1	0.01	0.23	0.19	0.03	
Topic2	0.21	0.07	0.48	0.02	
Topic3	0.53	0.01	0.17	0.04	

Each topic contains a score for all the words in the corpus.

Finding Representative Words for a Topic

- We can **sort the words** with respect to their probability score.

The top x words are chosen from each topic to represent the topic. If $x =$

10, we'll sort all the words in topic1 based on their score and take the top 10 words to represent the topic.

This step may not always be necessary because if the corpus is small we can store all the words in sorted by their score.

- Alternatively, we can **set a threshold** on the score. All the words in a topic having a score above the threshold can be stored as its representative, in order of their scores.





Photo by Anusha Barwa on Unsplash

Let's say we have 2 topics that can be classified as *CAT_related* and *DOG_related*. A topic has probabilities for each word, so words such as *milk*, *meow*, and *kitten*, will have a higher probability in the *CAT_related* topic than in the *DOG_related* one. The *DOG_related* topic, likewise, will have high probabilities for words such as *puppy*, *bark*, and *bone*.

If we have a document containing the following sentences:

“Dogs like to *chew* on *bones* and fetch sticks”.

“Puppies drink *milk*.”

“Both like to *bark*.”

We can easily say it belongs to topic *DOG_related* because it contains words *such as Dogs, bones, puppies, and bark*. Even though it contains the word *milk* which belongs to the topic *CAT_related*, the document belongs to *DOG_related* as more words match with it.

Assumptions:

- Each document is just a collection of words or a “bag of words”. Thus, the **order of the words** and the **grammatical role** of the words (subject, object, verbs, ...) are **not considered** in the model.
- Words like am/is/are/of/a/the/but/... don't carry any information about the “topics” and therefore can be eliminated from the documents as a preprocessing step. In fact, **we can eliminate words that occur in at least %80 ~ %90 of the documents**, without losing any information.

For example, if our corpus contains only medical documents, words like human, body, health, etc might be present in most of the documents and hence can be removed as they don't add any specific information which would make the document stand out.

- We **know beforehand how many topics** we want. ‘*k*’ is pre-decided.

- **All topic assignments except for the current word in question are correct**, and then updating the assignment of the current word using our model of how documents are generated

How does LDA work?

There are 2 parts in LDA:

- The *words that belong to a document*, that we already know.
- The *words that belong to a topic* or the probability of words belonging into a topic, that we need to calculate.

The Algorithm to find the latter

- Go through each document and randomly assign each word in the document to one of k topics (k is chosen beforehand).
- For each document d , go through each word w and compute :
 1. **$p(\text{topic } t \mid \text{document } d)$: the proportion of words in document d that are assigned to topic t .** Tries to capture how many words belong to the topic t for a given document d . Excluding the current word.
If a lot of words from d belongs to t , it is more probable that word w

belongs to t .

($\frac{\text{\#words in } d \text{ with } t + \alpha}{\text{\#words in } d \text{ with any topic} + k * \alpha}$)

2. **p(word w | topic t):** the proportion of assignments to topic t over all documents that come from this word w . Tries to capture how many documents are in topic t because of word w .

LDA represents documents as a mixture of topics. Similarly, a topic is a mixture of words. If a word has high probability of being in a topic, all the documents having w will be more strongly associated with t as well. Similarly, if w is not very probable to be in t , the documents which contain the w will be having very low probability of being in t , because rest of the words in d will belong to some other topic and hence d will have a higher probability for those topic. So even if w gets added to t , it won't be bringing many such documents to t .

- Update the probability for the word w belonging to topic t , as

$$p(\text{word } w \text{ with topic } t) = p(\text{topic } t \mid \text{document } d) * p(\text{word } w \mid \text{topic } t)$$

A layman's example

Suppose you have various photographs(*documents*) with captions(*words*). You want to display them in a gallery so you decide to categorize the photographs on various themes(*topics*) based on which you will create different sections in your gallery.



Photo by Soragrit Wongsu on Unsplash

You decide to create $k=2$ sections in your album — nature and city. Naturally, the classification isn't so clear as some photographs with city have trees and flowers while the nature ones might have some buildings in it. You, as a start, decide to assign the photographs which only have nature or city elements in them into their respective categories while you randomly assigned the rest.

You notice that a lot of photographs in nature have the word *tree* in their captions. So you concluded that the word *tree* and topic *nature* must be closely related.

Next, you pick the word *building* and check how many photographs are in *nature* because they have the word *building* in their caption. You don't find many and now are less sure about *building* belonging to the topic *nature* and associate it more strongly with the topic *city*.

You then pick a photograph which has the caption “**The tree is in front of the building and behind a car**” and see that it is in the category nature currently.

You then chose the word ***tree***, and calculate the first probability **$p(\text{topic } t \mid \text{document } d)$** : other words in the caption are *building* and *car*, most photographs having captions with *building* or *car* in it are in *city*, so you get

a low probability.

Now the second probability $p(\text{word } w | \text{topic } t)$: we know that a lot of photographs in nature have the word *trees* in it. So you get a high score here.

You update the probability of *tree* belonging in *nature* by multiplying the two. You get a lower value than before for *tree* in *topic* nature because now you have seen that *tree* and words such as *building/car* in the same caption, implying that trees can also be found in cities.

For the same reason, when you update the probability for *tree* belonging in *topic city*, you will notice that it will be greater than what it was before.

After multiple iterations over all the photographs and for each topic, you will have accurate scores for each word with respect to each topic. Your guesses will keep getting better and better because you'll conclude from the context that words such as buildings, sidewalk, subway appear together and hence must belong to the same topic, which we can easily guess is *city*.

Words such as mountains, fields, beach which might not appear together in a lot of captions but they do appear often without city words and hence will have higher scores for nature.

While words such as trees, flowers, dogs, sky will have almost the same probability of being in either as they occur in both topics.

As for the photograph, you see that it has 1 word (with average probability) from category nature and 2 words (with high probability) from city, you conclude, it belongs to city more strongly than it does to nature and hence you decide to add it in city.

Side note

The applications of LDA need not be restricted to Natural Language Processing. I recently implemented a paper where we use LDA(along with a Neural Networks) to extract the scene-specific context of an image. If you are interested in learning more about that please leave a comment or a message.

References

- If you are looking for an implementation of LDA, here is an amazing article on Towards Data Science which does it pretty well!
- See “A layman’s example” in Edwin Chen Blog to gain more intuition about LDA.

. . .

I'm glad you made it till the end of this article. 🎉

I hope your reading experience was as enriching as the one I had writing this.



Do check out my other articles here.

If you want to reach out to me, my medium of choice would be Twitter.

Sign up for The Daily Pick

By Towards Data Science

Hands-on real-world examples, research, tutorials, and cutting-edge techniques delivered Monday to Thursday. Make learning your daily ritual. [Take a look](#)

Your email



Get this newsletter

By signing up, you will create a Medium account if you don't already have one. Review our [Privacy Policy](#) for more information about our privacy practices.

Machine Learning

Naturallanguageprocessing

Topic Modeling

Lda

Probability

Discover Medium

Welcome to a place where words matter. On Medium, smart voices and original ideas take center stage - with no ads in sight. Watch

Make Medium yours

Follow all the topics you care about, and we'll deliver the best stories for you to your homepage and inbox. Explore

Become a member

Get unlimited access to the best stories on Medium — and support writers while you're at it. Just \$5/month. Upgrade

[About](#)

[Help](#)

[Legal](#)