# Topic Modeling for Scientific Documents

**Jethro Kuan**
WING-NUS

## Abstract

Topic models are statistical models, used to discover abstract topics that occur in a collection of documents. In this paper, we focus on the technique of topic modeling, particularly in its application to the domain of scientific documents. We discuss different approaches to topic modeling, and weigh in on their strengths and shortcomings.

## 1 Introduction

In this information age, the availability of knowledge is insufficiently met with the tools to navigate it. Tools like Semantic Scholar use machine learning to process scientific documents and extract meaningful structure, empowering researchers to discover papers more relevant to their work. One such technique is topic modeling.

## 2 Topic Modeling

Topic Modeling attempts to discover abstract topics within documents. Most topic models are generative models, assuming that latent variables govern the generative process of a document.

For example, a topic model trained on a corpora of scientific documents could learn topics such as "Neural Networks", "Biology" and "Medicine". Topics attribute a high probability to words that relate to the topic. For example, a "Neural Networks" topic would give attribute a high probability to "classifiers", and a low probability to "bacteria". The quotations around the topics indicate that the labels are human interpretations of what these topics may be. Automatic labeling of these topics are also areas of research. **CITE**

Topic modeling may discover that document in Figure 1 relates to "Medicine" and "Neural Networks", because these topics give rise to the colored words in the document.

In addition, we can view these abstract topics as a form of clustering. Researchers are better able to navigate the large corpora of scientific knowledge, by exploring documents that have similar topic proportions.

In this paper, we briefly discuss the predominant model: Latent Dirichlet Allocation (LDA) and its vari-
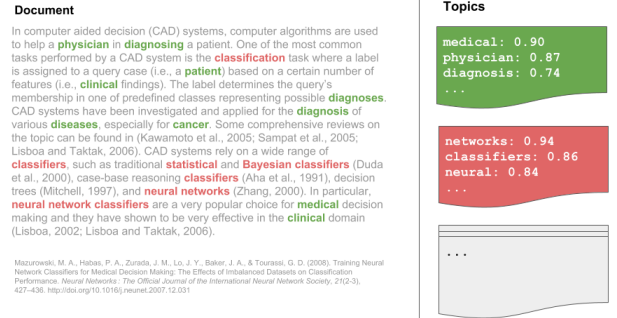


Figure 1: According to the trained topic model, topics like "Medicine" and "Neural Networks" generate the colored words of the document

ants. In addition, we explore other models that are not derivatives of LDA, such as the Replicated Softmax.

## 3 LDA

Latent Dirichlet analysis is widely considered to be the simplest topic model. LDA models each document as a mixture of topics, where a topic $\beta_k$ is a probability distribution over a fixed vocabulary of terms.

The Dirichlet distribution encourages sparsity, encoding the belief that the document-topic distribution has few topics per document, and the topic-word distribution has few words per topic. These two beliefs work against each other, and LDA discovers this sparsity balance, which gives rise to the structure of the textual data.

The generative process of a document is as follows:

```
for each document w do
  Draw topic distribution θ ∼ Dir(α);
  for each word at position n do
    Sample topic zₙ ∼ Multinomial(θ)
    Sample word wₙ ∼ Multinomial(βzₙ)
```

We can compute the probability of a word, given the LDA parameters:

$$p(w|\alpha,\beta) = \int_\theta \left( \sum_{n=1}^{N} \sum_{z_n=1}^{k} p(w_n|z_n,\beta)p(z_n|\theta) \right) p(\theta|\alpha)d\theta \tag{1}$$

Posterior inference over the hidden variables $\theta$ and $z$ is intractable due to the coupling between $\theta$ and $\beta$ under the multinomial assumption. (Blei et al., 2003). Hence, we rely on techniques for approximate inference of the posterior. We discuss these techniques in section 4.

## 3.1 Statistical Assumptions

As Mackay quips, we cannot make inference without statistical assumptions. LDA makes several assumptions, some rendering it less suited for application to the domain of scientific documents.

1. The order of documents do not matter.

The meaning of keyphrases used in scientific literature change over time. For example, neural networks meant a different thing two decades ago.

2. Bag of Words

Variants of LDA relax these statistical assumptions, or make other assumptions in place. We discuss Dynamic Topic Modeling (DTM) in section 5, and briefly mention the rest in the appendix.

## 4 Inference Methods

Approximating intractable probability densities is a well-studied problem in modern statistics. This problem arises often in Bayesian statistics, where computing posterior probablity densities in requires inference over latent variables. Many learning algorithms have been developed, including collapsed Gibbs Sampling, Variational Inference, Collapsed Variational Inference, and MAP estimation. Each of these approximation techniques have their own strength and shortcomings.

The two inference methods are briefly discussed below, and a comparison between them relegated to the appendix in subsection 6.1.

## 4.1 MCMC Sampling

Historically, Markov Chain Monte Carlo (MCMC) sampling has been the dominant technique for approximating posterior densities. In MCMC, we construct an ergodic Markov chain on the latent variable $z$, whose stationary distribution is the posterior $P(z|x)$. Samples are drawn from the stationary distribution, and used to approximate the posterior empirically.

## 4.2 Variational Inference

Mean field variational inference (MFVI) breaks the coupling between $\theta$ and $z$ by introducing free variational parameters $\gamma$ over $\theta$ and $\phi$ over $z$ and dropping the edges between them. This results in an approximate posterior $q(\theta,z|\gamma,\phi) = q_\gamma(\theta)\prod_n q_\phi(z_n)$.

To best approximate the true posterior, we frame it as an optimization problem, minimizing $L$ where:

$$L(\gamma,\phi|\alpha,\beta) = D_{KL}\left[q(\theta,z|\gamma,\phi)||p(\theta,z|\alpha,\beta)\right] - \log p(w|\alpha,\beta) \tag{2}$$

This optimization has closed form coordinate descent equations for LDA, because the Dirichlet is conjugate to the Multinomial distribution. This computational convenience comes at the expense of robustness, making it difficult to apply to other more complicated topic models.

## 5 Dynamic Topic Modeling

Dynamic Topic Modeling (DTM) was proposed to remove the assumption that documents are *exchangeable*. (Blei and Lafferty, 2006)

The order of documents are important for scientific documents, since both the content, and the meaning of words evolve over time.

In DTM, data is divided by discrete time slices. The topics associated with time slice $t$ evolve from the time slice $t-1$. Because the Dirichlet distribution is not amenable to sequential modeling, we use the Gaussian distribution to model the sequence of random variables.

The generative process for time slice $t$ is as follows:

```
Draw topic distribution βt|βt−1 ∼ N(βt−1, σ²I);
Draw αt|αt−1 ∼ N (αt−1, δ²I)
for each document w do
  Draw ηw,t ∼ N (αt, a²I)
  for each word at position n do:
    Sample topic zt,n ∼ Multinomial(π(ηw,t))
    Sample word wt,d,n ∼ Multinomial(βt,z,n)
```

$\pi$ maps the multinomial parameters to the mean parameters, $\pi\left(\beta_{k,t}\right)_w = \frac{exp(\beta_{k,t,w})}{\sum_w exp(\beta_{k,t,w})}$

The Multinomial and Guassian distributions are not conjugates, inference via Gibb's sampling is difficult. Hence, variational inference instead.

Further extensions of this approach include the continuous Dynamic Topic Models (cDTM), which removes the discretization of the time slices. (Wang et al., 2012) This model has been used to predict the timestamp of documents.

## 6 Appendix
### 6.1 Choosing an Inference Method

How do we know which technique to use to approximate the posterior density? MCMC methods are computationally more intensive, but provide samples that are approximately exact from the target posterior density. In contrast, VI methods view the problem as an optimization problem, which allows it to utilize efficient learning algorithms such as stochastic optimization. This is much quicker to compute, and is suited for larger datasets.

MCMC methods, however, cover a large family of sampling methods. Gibb's sampling requires that the prior and posterior are conjugate distributions. When this is not possible, such as in DTM, VI methods can

perform better than other methods in the MCMC family.

A closer look at the different inference approximation algorithms, however, shows that the performance differences can be explained away by setting certain smoothing hyperparameters (Asuncion et al., 2012).

## 6.2 Alternatives to LDA

# References

Arthur Asuncion, Max Welling, Padhraic Smyth, and Yee Whye Teh. On smoothing and inference for topic models. *CoRR*, 2012. URL `http://arxiv.org/abs/1205.2662v1`.

David M Blei and John D Lafferty. Dynamic topic models. In *Proceedings of the 23rd international conference on Machine learning*, pages 113–120. ACM, 2006.

David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.

Chong Wang, David Blei, and David Heckerman. Continuous time dynamic topic models. *CoRR*, 2012. URL `http://arxiv.org/abs/1206.3298v2`.