

Algorithmic Methods in the Humanities – Summer 2016

1 Latent Dirichlet Allocation

Florian Becker

Abstract

Due to the vast amounts of texts which are produced and digitized on a daily basis, it becomes more and more difficult to find what we are looking for. Topic models can be used as a tool to discover *topics* in large collections of texts. Each text is considered to be a distribution over topics and topics are, in turn, distributions over words. Topic model algorithms receive unlabeled texts as input and produce the topic distribution as output. Latent Dirichlet Allocation is a generative probabilistic topic model. The basic idea is that documents are mixtures of latent topics distributed according to a Dirichlet prior.

1.1 Introduction

When seeking digitized and hyper-linked information, we use searching tools that weigh each result according to some centrality measure. Such techniques are often sufficient, provided that we only want to find some result and do not care about the underlying topic structure which might give us more insight into a document or even a collection of many documents. Topic models enable us to not only search documents by the themes that they contain, but also give us a tool with which we can automatically organize and structure text corpora. It is assumed that a document exhibits multiple topics and each is a distribution over words. Consider, for instance, a corpus with texts about digital humanities. A topic model algorithm might output that the text corpus contains the topics *Computer Science*, *Linguistics* and *Philosophy*. But since every topic is a distribution over words, the algorithm will not just output '*Computer Science*', but will rather give a distribution over words, which we then can identify as the topic '*Computer Science*': e.g. *algorithm*, *program*, *complexity*, *machine*, *Turing*, *artificial intelligence*, etc.

Topic Models can also be thought of as a way of clustering similar documents. Scientific papers, books, tweets, blogs etc. can be organized according to the themes they feature. With topic models organizing a large corpus becomes feasible.

Latent Dirichlet Allocation (LDA) is a generative topic model. In topic models the topic and word distribution of a document can be explained by a generative model. In machine learning generative models are used to *simulate* data points. Given some data (e.g. a text corpus) in a generative model one will assume that there exists some process in the background which has generated the data. In our topic model context, this means that the process first decided what to pick as a topic and then decided which word to pick from that topic. The generative process has certain *parameters* and it is assumed that there must be a certain configuration that is most likely to have generated the observation. The goal is then to infer those hidden parameters. Hence, inference can be seen as reversing the generative process. That means that we want to uncover the most likely configuration that has produced a document. In the case of Latent Dirichlet Allocation the hidden variables are the per-document topic distribution, the per-document per-word topic assignments and the topics themselves.

Generally in *discriminative* machine learning data-driven algorithms are building a model.

This process is called *training* or *learning* and refers to the fact that the model adapts iteratively to the training data, such that eventually the model can explain the target variables. *Explaining* the data in a supervised model means, that the model delivers a label to a query. In unsupervised learning, a model represents the structure of unlabeled data. In a generative (probabilistic) model however, we understand something else under the term *model*. A *model* in this context means something that actively models the data/observations. Section 1.4 and 1.5 discuss the differences between discriminative and generative models in more detail.

Latent Dirichlet Allocation was introduced by David Blei, Andrew Ng, and Michael Jordan in 2003 [2]. Since then it has been applied in various different domains to many different sorts of texts. It is not restricted to operate on ‘classical’ documents. It was applied, for instance, in the domain of public health, to analyze tweets with LDA in order to track illnesses over time [8]. It was also just recently applied in bioinformatics for annotating genomes with a feature term that describes a feature with LDA [9] .

As the intersection of computing and humanities grows, topic models will become more important. Humanities profit from algorithmic methodologies applied on large text corpora for obvious reasons; manually structuring those is mostly not an option.

LDA can also be seen as a way of reducing the dimensionality of a corpus, since the topic distribution is a low dimensional representation of the documents [7]. Thus, one may gain insight about the semantic levels when examining the low-dimensional structure of a text corpus.

After a section about preliminaries and notation this report will introduce plate notation in section 1.3 and the generative model in section 1.5. Section 1.4 will discuss LDA in terms of machine learning and will therefore also introduce some terminology. In section 1.6 plate notation will be used in order to introduce a concise way to formulate Latent Dirichlet Allocation as a statistical model. Afterwards, the problem of Bayesian inference is discussed and a solution, the Gibbs Sampling algorithm, is presented.

1.2 Preliminaries and Notation

Throughout this report the notation of the original paper by Blei et al. will be used and is summarized in the following table.

■ **Table 1** Notation, as used in [2].

| symbol | description |
|---------------|--|
| α | parameter of the Dirichlet prior on the per-document topic distributions |
| β | parameter of the Dirichlet prior on the per-topic word distribution |
| θ_i | topic distribution for document i |
| φ_k | word distribution for topic k |
| z_{ij} | topic for the j -th word in document i , and |
| w_{ij} | specific word |
| \mathbf{w} | a sequence of words (a document) |
| \mathcal{D} | a corpus of M documents; $\mathcal{D} = \{w_1, \dots, w_M\}$ |

In order to understand in what way topics are assumed to be distributed the multinomial and Dirichlet distribution including important properties will be summarized briefly.

1.2.1 Multinomial Distribution

The binomial distribution is a discrete distribution over events with two outcomes. It gives the probability of obtaining k successes out of n trials, where each success has a probability of p . The probability mass function (PMF) is given by (Equation 1).

$$\Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k} \quad (1)$$

The binomial distribution is a special case of the multinomial distribution. The multinomial distribution models events which have k outcomes, where each has a fixed probability. The probability mass function is given by

$$f(x_1, \dots, x_k; n, p_1, \dots, p_k) = \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k} \quad (2)$$

1.2.2 Dirichlet Distribution

The Dirichlet distribution $\text{Dir}(\alpha)$ is a probability distribution parameterized by a positive real valued vector $\alpha = (\alpha_1, \dots, \alpha_K)$. The probability density function is given by:

$$f(x_1, \dots, x_K; \alpha_1, \dots, \alpha_K) = \frac{\Gamma\left(\sum_{i=1}^K \alpha_i\right)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K x_i^{\alpha_i-1} \quad (3)$$

where Γ is the Gamma function.

The Dirichlet distribution is often also expressed with the inverse of the Beta function as a normalizing constant. $B(\alpha)$ is constant as it only depends on the fixed parameter α .

$$f(x_1, \dots, x_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\alpha)} \prod_{i=1}^K x_i^{\alpha_i-1} \quad (4)$$

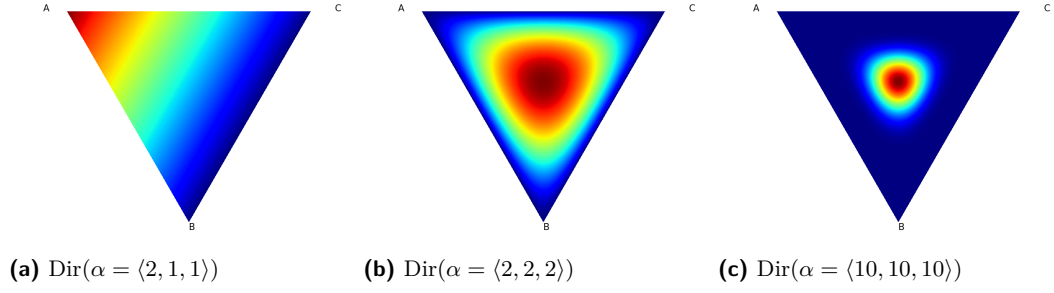
where

$$B(\alpha) = \frac{\prod_{i=1}^K \Gamma(\alpha_i)}{\Gamma\left(\sum_{i=1}^K \alpha_i\right)} \quad (5)$$

The expected value of the Dirichlet distribution is simply computed by:

$$\mathbb{E}[X_i] = \frac{\alpha_i}{\alpha_0}, \quad \alpha_0 = \sum_{i=1}^K \alpha_i. \quad (6)$$

The support of the Dirichlet distribution of order K is the $(K - 1)$ -simplex. For $K = 3$ the 2-simplex is given by a regular triangle, where each corner is an event. The vertices are $(0, 0, 1)$, $(0, 1, 0)$ and $(1, 0, 0)$. A point $\vec{p} = (p_1, p_2, p_3)$ on this probability simplex corresponds to a certain event.

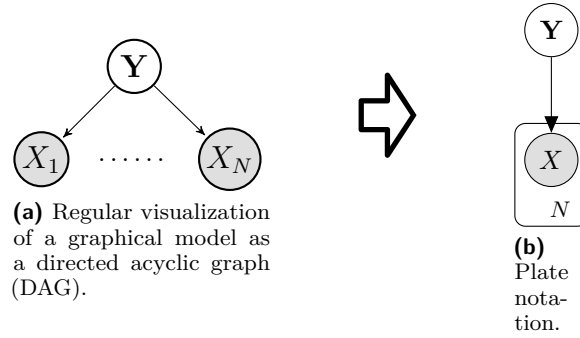


■ **Figure 1** Dirichlet distributions in the two-dimensional simplex with different α . The corners A , B , C correspond to topics. Regions of higher probability are shown by a heat map, i.e. (dark) red corresponds to the most probable sample.

1.3 Plate Notation

Generally, Bayesian networks (or probabilistic directed acyclic graphical models) represent conditional dependencies with a directed acyclic graph (DAG). A directed edge from a random variable Y to a random variable X means that Y *causes* X .

In a probabilistic graphical model conditional dependencies between various random variables are often depicted in plate notation. Plate notation is a concise method for visualizing the factorization of the joint probability. Typically, shaded nodes stand for *observed* random variables. The plates represent the repeated structure.



■ **Figure 2** Plate notation as a way to summarize conditional dependencies. Shaded nodes represent *observed* variables. Here, the X_i are conditionally dependent on Y . In plate notation all the X_i are summarized as one plate.

The joint probability of the probabilistic graphical model depicted in (Figure 2) is

$$\begin{aligned} \Pr[X_1, \dots, X_n] &= \prod_{i=1}^n \Pr[X_i | \text{parent}(X_i)] \\ &= \prod_{i=1}^n \Pr[X_i | Y] \end{aligned} \tag{7}$$

The plate notation will be used in the next section to have a concise representation of Latent Dirichlet Allocation and especially the conditional dependencies.

1.4 Machine Learning and Topic Models

In this section Latent Dirichlet Allocation will be discussed in terms of machine learning. Therefore, the distinction between supervised and unsupervised as well as the distinction between discriminative and generative model will be clarified.

Machine learning algorithms or problems are usually classified into two broad categories: supervised and unsupervised learning¹. In supervised learning the task is to find a hypothesis or model from labeled data such that the model can correctly classify the data. Given a set of *training examples* $\{x_{1:N}, y_{1:N}\}$, where $x_i \in X$ is an instance or *feature vector* and $y_i \in Y$ is its corresponding label a supervised machine learning algorithm must find a hypothesis $h : X \rightarrow Y$ such that $\forall i : h(x_i) = y_i$. In contrast to that, in an unsupervised machine learning problem there are no labels. Clustering, for instance, falls into the domain of unsupervised learning. Given data vectors $\{x_{1:N}\}$ the task is to label them. The underlying premise of clustering algorithms (such as *K-Means* for example) is that data that belongs to the same cluster is more *similar* to each other than to data from other clusters. Given an input space X and two clusters $V = \{v_{1:N} | v \in X\}$ and $W = \{w_{1:M} | w \in X\}$ and some notion of similarity $s : X \times X \rightarrow \mathbb{R}$, it must hold that $\forall v_1, v_2 \in V$ and $\forall w \in W : s(v_1, v_2) > s(v_1, w)$. The objective function is to maximize the similarities within one cluster and minimize the similarities between clusters. Latent Dirichlet Allocation is an unsupervised model, which can be seen as a clustering algorithm. The input space X consists of documents, which exhibit K topics. Unlike *K-Means*, in LDA a document can be part of more than one cluster (see also Figure 6).

1.4.1 Discriminative Models

Discriminative algorithms, as opposed to generative ones, directly optimize a model for a classification task [6]. Given an observed variable x and unobserved variable y , discriminative models will give the conditional distribution $P(y|x)$. *K-Means* for instance is a similarity-bases, unsupervised, discriminative approach [11]. Examples of discriminative models include *linear regression*, *support vector machines* and *nerual networks*. However, LDA is a generative model. The assumption is that a corpus is generated by a process, which will be discussed in the next section.

1.5 Generative Process

Generative models generate data values according to a probability distribution. In machine learning generative models are used to estimate the joint probability distribution $Pr(X, Y)$ of observed data X and labels Y . In the context of probabilistic topic models, Y corresponds to the latent variables (i.e. to the unknown parameters) and X to the (observable) words in a corpus. Latent Dirichlet Allocation assumes the following generative process for a corpus:

¹ Reinforcement learning does not fall into this categorization. Since it is not of importance in this context, it won't be considered here.

Input : Number of words N , Number of topics K
 Dirichlet parameters $\alpha \in \mathbb{R}_+^K, \beta \in \mathbb{R}_+^N$
Output : D documents (bag-of-words) with N words each

1. Choose $\theta_i \sim \text{Dir}(\alpha)$,
2. Choose $\varphi_k \sim \text{Dir}(\beta)$

for each word position i, j **do**
 (a) Choose a topic $z_{i,j} \sim \text{Multinomial}(\theta_i)$.
 (b) Choose a word $w_{i,j} \sim \text{Multinomial}(\varphi_{z_{i,j}})$
end

Algorithm 1 : Generative Model

The generative model contains the main ideas of LDA. Data is treated as observations coming from a process which samples words from a (hidden) distribution over topics. Every word is generated independently from any other word. The generative process does not take word order into account. Hence, the generative model produces a bag-of-words, where word order does not play a role.

If three topics are chosen, e.g. *Linguistics*, *Philosophy* and *Computer Science*, and if equal probability is given to all of them, then the generative process might produce bag-of-words resembling a text from *Digital Humanities*. By putting more weight on *Linguistics* and *Computer Science* the result would resemble documents from the domain of *Computer Linguistics*. The generative process of LDA is depicted as a Bayesian probabilistic model in Figure 4.

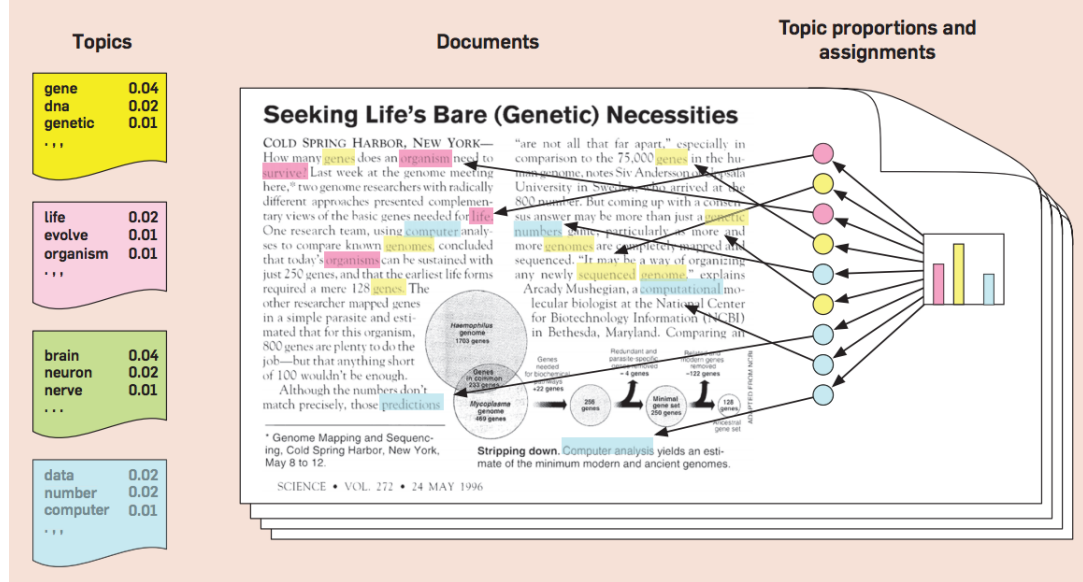
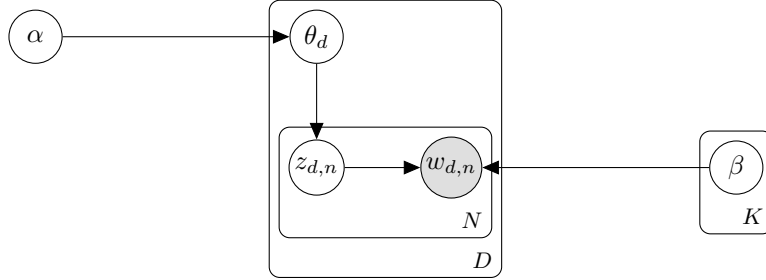


Figure 3 The generative process. Every word originates from a topic coming from a distribution over topics. The topics and their word distributions are depicted on the far left. E.g. the *yellow* topic is made up of the word *gene* by 4%. Taken from [1]

In Figure 3 the generative model is illustrated. The topic proportions are plotted as a

bar chart, the colored coins are the topic assignments for each word.

1.6 LDA as a graphical model



■ **Figure 4** Plate notation for Latent Dirichlet Allocation corresponding to the generative process. Every vertex depicts a random variable. The shaded node $w_{d,n}$ stands for the observed words. Everything else is not observed and must be inferred.

From the Bayesian probabilistic model one can derive the probability of a corpus $\mathcal{D} = \{\mathbf{w}_1, \dots, \mathbf{w}_n\}$. First, the joint distribution of a mixtures of topics θ , N topics \mathbf{z} and N words \mathbf{w} can be calculated by considering the conditional structure given in the graphical model:

$$p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) = p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n, \beta) \quad (8)$$

In order to arrive at the marginal distribution of one single document \mathbf{w} , one must integrate over the (continuous) random variable θ and sum over the topic variable z .

$$p(\mathbf{w} | \alpha, \beta) = \int p(\theta | \alpha) \left(\prod_{n=1}^N \sum_{z_n} p(z_n | \theta) p(w_n | z_n, \beta) \right) d\theta \quad (9)$$

The probability of the corpus $p(\mathcal{D} | \alpha, \beta)$ is now obtained by multiplying all the marginal probabilities of the M documents:

$$p(\mathcal{D} | \alpha, \beta) = \prod_{d=1}^M \int p(\theta_d | \alpha) \left(\prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn} | \theta_d) p(w_n | z_{dn}, \beta) \right) d\theta_d \quad (10)$$

1.7 (Bayesian) Inference

In Bayesian statistics, one is interested in computing the degree of belief $p(h | D)$ of a hypothesis $h \in \mathcal{H}$ given data D , where \mathcal{H} is the space of all hypotheses. Bayes theorem states:

$$p(h | D) = \frac{p(h)p(D|h)}{p(D)} \quad (11)$$

$p(h | D)$ is also called the posterior.

Inference means computing the *posterior* and corresponds to ‘*reversing*’ the generative model. In other words, inference is about fitting a generative model (or rather the hidden

parameters), such that it explains the observed data. In Bayesian statistics, maximum likelihood estimation is used in order to estimate the parameters of a model given data. Let τ be the parameter of a statistical model and let $\mathbf{x} = x_1, \dots, x_n$ be the observations, then the maximum likelihood estimator is defined to be [3]:

$$\arg \max_{\tau} \mathcal{L}(\tau; \mathbf{x}) = \arg \max_{\tau} P(\mathbf{x}|\tau) \quad (12)$$

Many methods solving the problem of maximum a posteriori estimation have been proposed. The original paper on LDA [2] suggested a Variational Bayesian method, which directly maximizes $P(\mathbf{w}|\varphi, \theta)$, which corresponds to equation 12 and can be understood as an extension of the expectation-maximization algorithm [4]. Gibbs Sampling as a way to compute the posterior was proposed by Steyvers and Griffiths [10]. They considered the posterior distribution over the assignments of words $w_{1:D}$ to topics $z_{1:D}$:

$$p(\beta_{1:K}, \theta_{1:D}, z_{1:D} | w_{1:D}) = \frac{p(\beta_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D})}{p(w_{1:D})} \quad (13)$$

The posterior is intractable to compute. More precisely, the denominator of (Equation 13), the marginal probability, is not feasible to compute. To compute $p(w_{1:D})$ one would have to sum up the joint distribution over all possible instances of the hidden topic structure.

1.7.1 Gibbs Sampling

Gibbs Sampling, named after the physicist Josiah Willard Gibbs, is a Markov Chain Monte Carlo Algorithm (MCMC) and can be used to approximate the marginal and posterior distribution by constructing a Markov chain which converges to the target distribution. A Markov chain represents a random process, where going from one state to another is determined by a transition matrix P . An entry $p_{i,j}$ corresponds to the probability of going from state i to state j . The goal of Gibbs Sampling is to integrate out the per-document topic proportions θ . In order to assign a word to a topic, the Gibbs sampler computes the probability of a topic $z_{d,n}$ if it is assigned to a word $w_{d,n}$, where all other word-to-topic assignments are given or fixed. $z_{-d,n}$ is defined to be the notation for all topic assignments except for for $z_{d,n}$

Formally:

$$p(z_{d,n} = k | z_{-d,n}, w, \alpha, \beta) = \frac{p(z_{d,n} = k, z_{-d,n} | w, \alpha, \beta)}{p(z_{-d,n} | w, \alpha, \beta)} \quad (14)$$

It was shown [5] that Equation 14 can be calculated by:

$$p(z_{d,n} = k | z_{-d,n}, w, \alpha, \beta) \propto \frac{C_{n,k}^{WT} + \beta}{\sum_{n=1}^W C_{n,k}^{WT} + W\beta} \cdot \frac{C_{d,k}^{DT} + \alpha}{\sum_{t=1}^T C_{d,t}^{DT} + T\alpha} \quad (15)$$

where T is the number of topics, and W the number of unique words. $C_{n,k}^{WT}$ and $C_{d,k}^{DT}$ are matrices. $C_{n,k}^{WT}$ is the number of times topic k was assigned to word n . On the other hand $C_{d,k}^{DT}$ is the number of times topic k was assigned to words in document d . The first step of the Gibbs Sampling algorithm is to randomly assign every word of all documents to a topic. By this the count matrices $C_{n,k}^{WT}$ and $C_{d,k}^{DT}$ are filled with values. Gibbs Sampling will

now sample topic assignments according to Equation 15. The first term of the product in Equation 15 also corresponds to an estimate of φ_i , and the second term to an estimate of θ_j

$$\hat{\varphi}_i = \frac{C_{i,j}^{WT} + \beta}{\sum_{k=1}^W C_{k,j}^{WT} + W\beta} \quad (16)$$

$$\hat{\theta}_j = \frac{C_{d,j}^{DT} + \alpha}{\sum_{t=1}^T C_{d,t}^{DT} + T\alpha} \quad (17)$$

► **Example 1.** The following step-by-step example illustrates how Gibbs Sampling can be used to assign a topic to a word. Consider the sentence (which represents one document out of a corpus):

Text mining algorithms can be used to find structure in text corpora like Plato's dialogues

A priori we choose the number of topics $K = 3$. The first step is to remove stop words and randomly assign each word to a topic (Table 2). This is done to all documents.

■ **Table 2** random initialization of topic assignments to words.

| | | | | | | |
|------|--------|------------|-----------|---------|-------|-----------|
| 1 | 3 | 2 | 1 | 2 | 1 | 2 |
| text | mining | algorithms | structure | corpora | Plato | dialogues |

The counts for all documents is shown in Table 3.

■ **Table 3** word counts for every topic

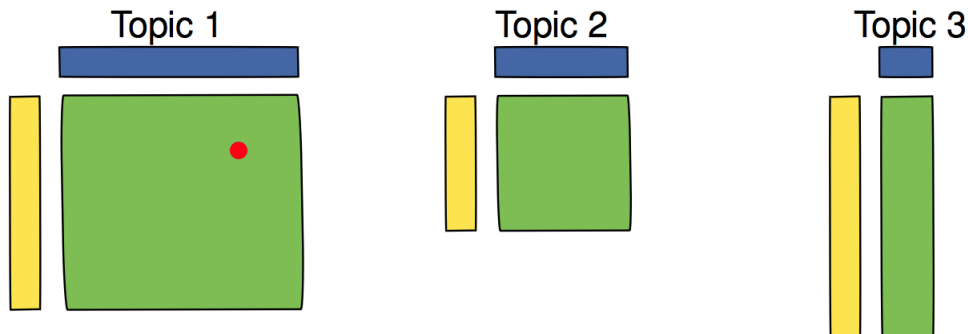
| | 1 | 2 | 3 |
|------------|-----|----|-----|
| text | 65 | 54 | 59 |
| mining | 21 | 4 | 12 |
| algorithms | 100 | 74 | 122 |
| structure | 20 | 12 | 14 |
| corpora | 5 | 2 | 12 |
| Plato | 35 | 33 | 42 |
| dialogues | 24 | 27 | 31 |

The goal is now to resample the word *algorithm*.

■ **Table 4** We first choose *algorithm* to be resampled.

| | | | | | | |
|------|--------|------------|-----------|---------|-------|-----------|
| 1 | 3 | ??? | 1 | 2 | 1 | 2 |
| text | mining | algorithms | structure | corpora | Plato | dialogues |

The assignment of the word *algorithm* is determined by sampling according to the topic distribution in the document and the word distribution over the corpus. In Figure 5 the sampling is depicted. The blue bars correspond to the topic distribution, the yellow bars to the word distribution over topics. Accordingly, it is most likely that *algorithm* is assigned to topic 1, since the green area for topic 1 is the largest. After assigning the word, the word counts for each topic get updated and the first iteration of the Gibbs sampling algorithm is completed. This process is repeated several times, until the assignments do not change anymore.



■ **Figure 5** Sampling according to the green area. Blue bars represent the topic distribution, the yellow bars to the word distribution over topics.

1.8 Experiments and Applications

To fully grasp what LDA will produce as output, this section will provide an example with *real world* data.

Gensim² is a tool for unsupervised semantic modeling. The idea is to define three topics by using Wikipedia articles. Then, all those Wikipedia articles are put into one corpus to see whether LDA can reproduce those three original topics. Basically this is a clustering task with a predefined number of clusters K . Similar to *K-Means* LDA will use the training set (the predefined Wikipedia articles) and will find the clusters that maximize the model parameters.

With Gensim the model can be trained very easily³:

■ Listing 1 Training the LDA model

$K = 3$

² <https://radimrehurek.com/gensim/index.html>

³ the complete ipython notebook can be found here: https://github.com/flobeck/algo-seminar/blob/master/ipython-notebook/lda_mixture.ipynb

```
lda = ldamodel.LdaModel(CORPUS, id2word=dictionary, num_topics=K,
update_every=1, chunksize=1, passes=350)
```

Gensim uses the Gibbs sampling procedure that was discussed in section 1.7.1. CORPUS contains the following preprocessed Wikipedia articles. The preprocessing steps comprise removing stop words (i.e. most common words) and transforming the documents to a bag-of-words.

- *Computer science* Algorithm, Computation, Computer Programming, Programming Language, Computational complexity theory, Computability theory, Artificial Intelligence
- *Philosophy* Epistemology, Metaphysics, Continental philosophy, Ancient Greek, Ethics, Aesthetics, Art, Phenomenology (philosophy), Pythagoras, Plato
- *Linguistics* Language, Semantics, Syntax, Phonology, Grammar, Phonetics, Pragmatics, Corpus linguistics, Linguistic prescription, Linguistic description

The output of a call of LDA with $K=3$ is summarized in table 5.

Given a set corpus of Wikipedia articles, LDA finds three topics which resemble the topics *computer science*, *philosophy* and *linguistics*.

Furthermore, it is possible to query where on the topic simplex an unseen document is located:

■ Listing 2 Query for a document

```
noam = wikipedia.page("Noam Chomsky").content
bow_vector = dictionary.doc2bow(tokenize(noam))
print lda[bow_vector]

>> [(1, 0.032599745516170064),
(2, 0.30504589551791855),
(3, 0.66235435896591133)]
```

In this case (Listing 2), the Wikipedia article on the linguist *Noam Chomsky* is made up of 66% Topic 3 (*Linguistics*), 30% Topic 2 (*Computer Science*) and 3% Topic 1 (*Philosophy*). Blei et. al showed that one can train a model for document classification with Latent Dirichlet Allocation [2]. Instead of using individual words as features, one can use the topic proportions for each document. A document can then be classified by using LDA to extract the topic proportion feature vectors of the training set and then a finding decision boundary with a Support Vector Machine (SVM) or some other linear classifier.

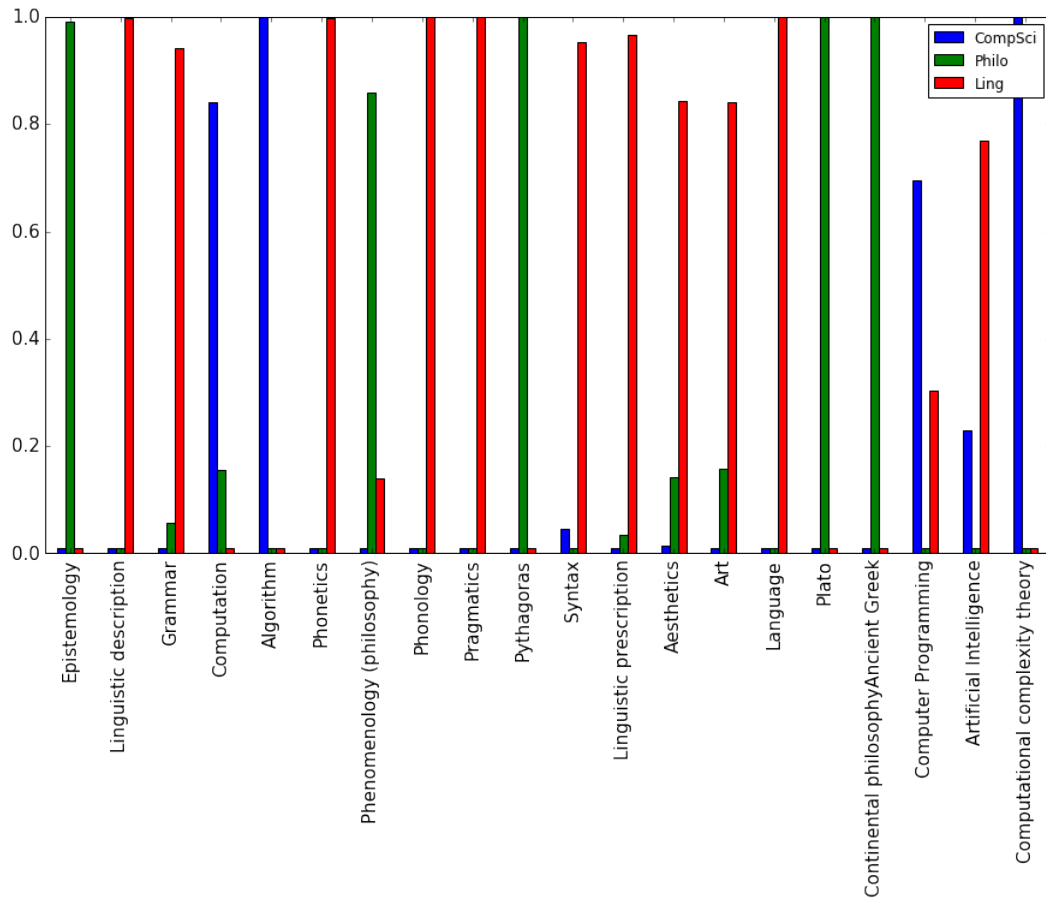


Figure 6 Topic proportions for some of the Wikipedia articles that were used to train the model. The Wikipedia article on ‘Pythagoras’, for instance, is classified as a document in the category ‘philosophy’ with very high confidence. Note that albeit the confidence is very high it is not 100%. Due to the model every document exhibits all the K topics, but to a very different degree.

■ **Table 5** Each of the three topics is a distribution over the words. E.g. topic 1 resembles *computer science*. The table only shows the most relevant words for each topic.

| TOPIC | REL. FREQUENCIES | WORD |
|-------|------------------|---------------|
| 1 | 0.010 | programming |
| | 0.009 | turing |
| | 0.008 | problem |
| | 0.007 | machine |
| | 0.006 | problems |
| | 0.006 | algorithms |
| | 0.006 | algorithm |
| | 0.006 | set |
| | 0.006 | complexity |
| | 0.005 | time |
| 2 | 0.008 | plato |
| | 0.006 | plato's |
| | 0.006 | knowledge |
| | 0.005 | pythagoras |
| | 0.005 | philosophy |
| | 0.004 | phenomenology |
| | 0.004 | greek |
| | 0.004 | according |
| | 0.004 | socrates |
| | 0.003 | theory |
| 3 | 0.019 | language |
| | 0.009 | languages |
| | 0.007 | art |
| | 0.005 | meaning |
| | 0.005 | linguistic |
| | 0.005 | human |
| | 0.004 | ethics |
| | 0.004 | speech |
| | 0.004 | study |
| | 0.004 | grammar |

1.9 Conclusion

Topic models enable scholars to find the hidden topic structure in large text corpora. The hidden structure must be uncovered in order to gain information about the low-dimensional structure behind a text corpus. A topic modeling algorithm takes a corpus as an input and will output the topics running through that corpus. The generative process is the essential idea behind Latent Dirichlet Allocation. Words come from a distribution over topics, which are again distribution over words. Reversing the generative model by applying Gibbs sampling yields a solution of the inference problem, which can not be solved analytically.

References

- 1 David M Blei. Probabilistic topic models. *Communications of the ACM*, 55(4):77–84, 2012.
- 2 David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- 3 Richard O Duda, Peter E Hart, and David G Stork. *Pattern classification*. John Wiley & Sons, 2012.
- 4 Charles W Fox and Stephen J Roberts. A tutorial on variational bayesian inference. *Artificial intelligence review*, 38(2):85–95, 2012.
- 5 Thomas L Griffiths and Mark Steyvers. Finding scientific topics. *Proceedings of the National academy of Sciences*, 101(suppl 1):5228–5235, 2004.
- 6 Tony Jebara. *Discriminative, generative and imitative learning*. PhD thesis, Massachusetts Institute of Technology, 2001.
- 7 Tomonari Masada, Senya Kiyasu, and Sueharu Miyahara. Comparing lda with plsi as a dimensionality reduction method in document clustering. In *Large-Scale Knowledge Resources. Construction and Application*, pages 13–26. Springer, 2008.
- 8 Michael J Paul and Mark Dredze. You are what you tweet: Analyzing twitter for public health. *ICWSM*, 20:265–272, 2011.
- 9 Pietro Pinoli, Davide Chicco, and Marco Masseroli. Latent dirichlet allocation based on gibbs sampling for gene function prediction. In *Computational Intelligence in Bioinformatics and Computational Biology, 2014 IEEE Conference on*, pages 1–8. IEEE, 2014.
- 10 Mark Steyvers and Tom Griffiths. Probabilistic topic models. *Handbook of latent semantic analysis*, 427(7):424–440, 2007.
- 11 Shi Zhong and Joydeep Ghosh. Generative model-based document clustering: a comparative study. *Knowledge and Information Systems*, 8(3):374–384, 2005.