

1 Discussed models

This section will introduce the reader into the two Topic modeling approaches which will be compared in this Thesis. The aim of both procedures is to assign one or more topics to different documents. Even if the vocabulary and the notation are similar for both approaches, the notation should be resumed at the beginning of the description of the respective model. The basic structural notation of the data consists of the following variables.

A collection of documents is called corpus $D = (\mathbf{w}_1, \dots, \mathbf{w}_M)$. It consists out of M documents $\mathbf{w} = (w_1, \dots, w_N)$ which are itself separated in words w_i . These words are vectors of length V . V refers to the length of a vocabulary which holds all the words occurring in the corpus. The vector for a specific word w_i contains all 0 except for index $j \in \{1, \dots, V\}$ which represents this very one word in the vocabulary. This notation may indeed be extended through the addition of indices for documents, but this is not done here or in the standard literature on topic models due to its unnecessary complexity.

1.1 LDA model

Latent Dirichlet Allocation is a Bayesian approach and is often associated with the class of hierarchical models [A. Gelman, 2014]. The idea is based on the representation of exchangeable random variables (acc. to de Finetti) as mixture of distributions. Given that documents \mathbf{w} and words w_i in each document - both considered as random variables in this setting - are exchangeable in such a way, a mixed model such as the LDA model is appropriate [Blei, 2003].

The following notation is used in conjunction with the LDA model. Let z_j be the topics with $j \in \{1, \dots, k\}$. In the LDA setting we assume for every topic z_j there is a term distribution

$$\beta_j \sim Dir(\delta)$$

We further assume each document \mathbf{w} has a distribution of topics.

$$\theta \sim Dir(\alpha)$$

Then each word w_i of \mathbf{w} is generated by the following process:

1. Choose $z_i \sim Mult(\theta)$
2. Choose $w_i \sim Mult(\beta_{z_i})$ This distribution will be referred to as $p(w_i|z_i, \beta)$

You can summarize this setup in a plate diagram as shown in figure 1. The notation above, which is also used within the diagram, coincides with the

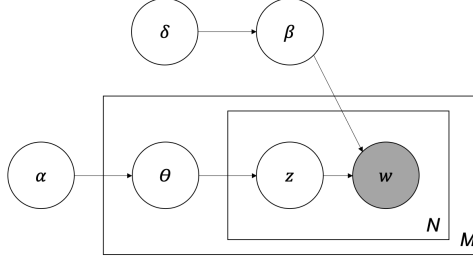


Figure 1: The well-established plate diagram for the standard LDA model extended by the parameter δ . The slightly bigger box represents the generative model of the corporis M documents. The smaller plate represents the iterative generation process of the N words of each document with the aid of the topics. See also "smoothed LDA model" in [Blei, 2003] for comparisons.

notation of [Hornik, 2011].

In order to estimate the model parameters, the first task is to calculate the posterior distribution, which consists of the joint distribution in the numerator and the marginal distribution in the denominator.

$$p(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta) = \frac{p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta)}{p(\mathbf{w} | \alpha, \beta)} \quad (1)$$

The joint distribution numerator can be derived straight forward.

$$p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) = p(\theta | \alpha) \prod_{i=1}^N p(w_i | z_i, \beta) p(z_i | \theta) \quad (2)$$

One can obtain the marginal distribution of a document \mathbf{w} , by integrating out the parameter θ and summing over the topics z_j . Nevertheless, this expression is intractable.

$$p(\mathbf{w} | \alpha, \beta) = \int p(\theta | \alpha) \left(\prod_{i=1}^N \sum_{z_i} p(z_i | \theta) p(w_n | z_i, \beta) \right) d\theta \quad (3)$$

The literature divides the approaches to calculating posterior distribution into two main categories. [Blei, 2012] distinguishes between sampling based algorithms and variational algorithms. [Powieser, 2012] lists a total of 6 algorithms that can be used to estimate parameters in the LDA model. This thesis will be confined to the two most cited and most used members of the two main groups. One approach is to simulate the posterior density

by iteratively sampling - the so-called Gibbs sampling. The second approach is a deterministic method, a modified version of the well-known EM algorithm [AP Dempster, 1977]: the Variational EM algorithm (VEM algorithm) [Wainwright and Jordan, 2008]. In the following two sections the both approaches are roughly outlined to give the reader some insight into the Bayesian inference underlying the algorithms.

1.1.1 Variational EM algorithm

In the VEM algorithm for the LDA model is a mean field approach which varies the steps E and M of the EM algorithm in a way such that this algorithm becomes solvable. Note that the main problem of calculating the marginal distribution is, to derive the conditional probability of some hidden variables given some observed values ("evidence"). The variation of the EM algorithms consists mainly in approximating the directly intractable E step. Rewriting the log of the border density of \mathbf{w} as follows in (4), results in the fact that the marginal density can be estimated downwards with the aid of Jensen's inequality.

$$\log p(\mathbf{w}|\alpha, \beta) = \log \int \sum_{\mathbf{z}} p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta) d\theta \quad (4)$$

$$= \log \int \sum_{\mathbf{z}} \frac{p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta) q(\theta, \mathbf{z})}{q(\theta, \mathbf{z})} d\theta \quad (5)$$

$$\geq \int \sum_{\mathbf{z}} q(\theta, \mathbf{z}) \log p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta) d\theta - \int \sum_{\mathbf{z}} q(\theta, \mathbf{z}) \log q(\theta, \mathbf{z}) d\theta \quad (6)$$

$$= \mathbb{E}_q[\log p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta)] - \mathbb{E}_q[\log q(\theta, \mathbf{z})] \quad (7)$$

Here $q(\theta, \mathbf{z})$ is an arbitrary distribution which can be called the variational distribution.

$$q(\theta, \mathbf{z}) \triangleq q(\theta, \mathbf{z}|\gamma, \phi) = q(\theta|\gamma) \prod_{i=1}^N q(z_i|\phi_i) \quad (8)$$

The right hand side $L(\gamma, \phi, \alpha, \beta) := \mathbb{E}_q[\log p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta)] - \mathbb{E}_q[\log q(\theta, \mathbf{z})]$ be called "lower bound". It can be shown that $\log p(\mathbf{w}|\alpha, \beta) - L(\gamma, \phi, \alpha, \beta)$ is the Kullbak Leibler divergence (D_{KL}) of the true posterior and the variational distribution. From equations (4)-(7) follows that:

$$\log p(\mathbf{w}|\alpha, \beta) = D_{KL}(q(\theta, \mathbf{z}|\gamma, \phi) || p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta)) + L(\gamma, \phi, \alpha, \beta) \quad (9)$$

Since the marginal is fixed, we conclude, that minimizing the KL-divergence is equivalent to maximizing the lower bound (see [Jordan, 1999] and [Wainwright and Jordan, 2008]),

for details of the derivation of the lower bound see [Blei, 2003]).

$$(\gamma^*, \phi^*) = \underset{\gamma, \phi}{\operatorname{argmin}} D_{KL}(q(\theta, \mathbf{z}|\gamma, \phi)||p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta)) \quad (10)$$

$$= \underset{\gamma, \phi}{\operatorname{argmax}} L(\gamma, \phi, \alpha, \beta) \quad (11)$$

The variation of the EM algorithm thus is to use the variational distribution $q(\theta, \mathbf{z}|\gamma^*(\mathbf{w}), \phi^*(\mathbf{w}))$ instead the posterior distribution $p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta)$. Now the two steps of the VEM algorithm are:

- (1) **E step** Optimize the variational parameters θ and ϕ for every document in the corpus. This can be done analytically by deriving the derivatives of the KL divergence. And set them to zero.
- (2) **M Step** Maximize the lower bound using the optimized parameter of the E step with respect to α and β .

1.1.2 Gibbs sampling

The second method to approximate the posterior distribution is Gibbs sampling, a so-called monte Carlo method. Instead of calculating the distributions for β and θ , the primary task is to find the posterior distribution over \mathbf{z} given the document \mathbf{w} . Gibbs sampling is also known as a Markov Chain Monte Carlo method. The name refers to the simulation process by which a chain of values is simulated whose limiting distribution desirably converges against the true distribution [M. Steyvers, 2006]. (12) shows the distribution, which is sampled from iteratively.

$$p(z_i = j | z_{-i}, w) \propto \frac{n_{-i,j}^{(l)} + \delta}{\sum_t n_{-i,j}^{(t)} + V\delta} \frac{n_{-i,j}^{(d_i)+\alpha}}{n_{-i}^{(d_i)} + k\alpha} \quad (12)$$

$z_i = j$... word-topic assignment of word i to topic j
 z_{-i} ... vector of word-topic assignments without the entry for word i
 $n_{-i,j}^{(l)}$... number of times the l th word in the vocabulary is assigned to topic j , not including the assignment for word i
 d_i ... document in the corpus which includes word i
 δ, α ... parameters of the prior distributions for β and θ

Usually the word-topic distributions $\beta_j^{(l)}$ for the words $l = 1, \dots, V$ and topics $j = 1, \dots, k$ and topic-document distributions $\theta_j^{(d)}$ for the documents $d = 1, \dots, D$ and the topics $j = 1, \dots, k$ will be of interest. (13) and (14)

shows the predictive distributions denoted as "estimators".

$$\hat{\beta}_j^{(l)} = \frac{n_{-i,j}^{(l)} + \delta}{\sum_t n_{-i,j}^{(t)} + V\delta} \quad (13)$$

$$\hat{\theta}_j^{(d)} = \frac{n_{-i,j}^{(d_i)+\alpha}}{n_{-i}^{(d_i)} + k\alpha} \quad (14)$$

For derivation and more details regarding the Gibbs sampling procedure see [M. Steyvers, 2006].

1.1.3 Implementation

In this thesis, the implementation of the LDA model and its estimation is mainly based on using the package `topicmodels` of Kurt Hornik. The package `topicmodels` can apply both the VEM algorithm as well as Gibbs sampling in order to fit the model. In addition, the package `tidytext` is being used for text structuring and embedding. Whereby there are other packages besides this implementation of the LDA model, `topicmodels` is particularly convenient, because `tidytext` was designed by its developers to work perfectly in combination with `topicmodels` [Silge and Robinson, 2017, p. 89].

1.2 Artificial Neural Networks

Artificial neural networks (ANN) are much more versatile than the LDA model. There are not only various forms of artificial neural networks, but also a very large number of application areas. Quite as machine learning procedures in general, also deep learning algorithms are divided into two broad categories: supervised learning, where a superset instance provides the algorithm with the output required to learn, and unsupervised procedures that internally train predefined models to find patterns in the input signals. In this chapter we will focus heavily on the former group of ANNs. Also, this chapter is intended to give the reader an overview of the research on neural networks as well as the background of their development.

1.2.1 Development

Research on ANNs dates back to the 1940s, when [McCulloch and Pitts, 1943] introduced the so called "M-P neuron". Whereby this neuron had only a bi-variate input and output, Rosenblatt later extended this idea to a network of M-P neurons, which allowed to set up a simple classification algorithm [Rosenblatt, 1958]. A perceptron in its basic form (single perceptron) is a binary classifier.

Imagine input data of a simple perceptron in the form of a matrix.

$$X = \begin{bmatrix} x_{11} & \dots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nn} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}$$

The dependent variable thus is a vector $\mathbf{y} = y_1, \dots, y_n$, with $y_i \in \{0, 1\}$. Consider the lines of the X matrix as vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$, with $\mathbf{x}_i \in \mathbb{R}^k$. The entries of each of the vectors are weighted with $\mathbf{w} = w_1, \dots, w_k$ with $w_j \in \mathbb{R}$ and aggregated in a function h e.g. a sum.

$$h(\mathbf{x}, \mathbf{w}) = \sum_{j=1}^k x_j w_j$$

Using a so called "activation function" h is mapped to the output space, which is in this case $O = \{0, 1\}$. At this point a step function serves as activation function.

$$a \circ h(\mathbf{x}, \mathbf{w}) = \begin{cases} 0 & \text{if } h(\mathbf{x}, \mathbf{w}) \leq 0 \\ 1 & \text{else} \end{cases}$$

The matrix X is passed vector by vector to the perceptron and the output is compared with the values for \mathbf{y} . During this procedure the weights are iteratively tuned by a simple updating algorithm using the pairs \mathbf{x}_i and y_i .

The algorithm of the simple perceptron is schematically shown in Figure 2. This diagram corresponds to the common representation in education [Mukherjee, 2019], although a horizontal perspective is often chosen.

If this basic perceptron is used in a clever way, designs can be developed that have variable application possibilities with this approach only. For example, a sigmoid function can be used as an activation function instead of the step function. So the output layer will not project into the $\{0, 1\}$ space, but into a probability space. If you add several nodes with the sigmoid activation function rather than a single output, you basically obtain the architecture that is also called multivariate logistic regression [Bahjat, 2006]. In that case the model can be estimated by an individual calculation of logistic regressions for each node in the output layer. This allows to classify not only two classes, but an arbitrary number (see the network architecture of Figure 3).¹

Shortly following the publication of these results, which would lay the foundation for later neural networks, ANN research had to suffer a severe setback

¹Note that in this setting the dependent variable must be one-hot encoded.

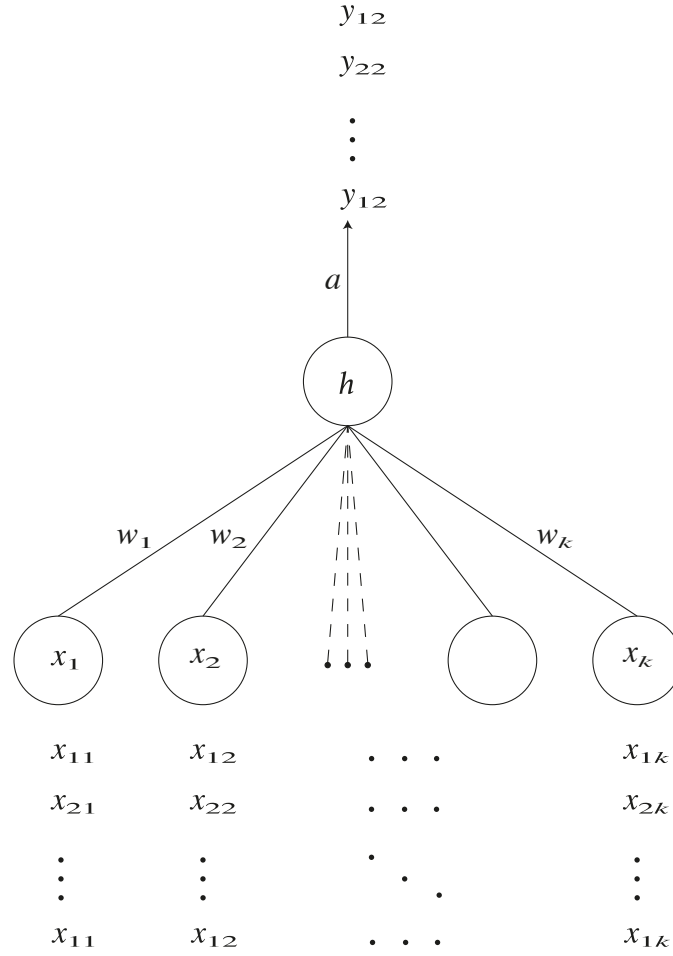


Figure 2: Schematic diagram of a simple perceptron by [Rosenblatt, 1958]

after Minsky and Papert [Minsky and Papert, 1969] were able to prove that perceptrons cannot provide a suitable solution in certain basic scenarios. It is possible to separate two clusters by a hyperplane with the perceptron, however, if the two groups could not be separated completely, the perceptron failed. For instance, it was impossible to find a solution for data generated with an *x-or* function (also called "exclusive or" function).

Although the researchers Minsky and Papert showed that the perceptron in this form cannot solve the "xor problem", they argued that extending the simple perceptron to a multi layer perceptron is an approach to solve this problem, if it was feasible to train this model. Instead of just one input layer and one output layer, an MLP may include an arbitrary number of hidden layers in between. Figure 4 shows the structure of such a network. However, with the naive optimizing algorithm of the perceptron as it was

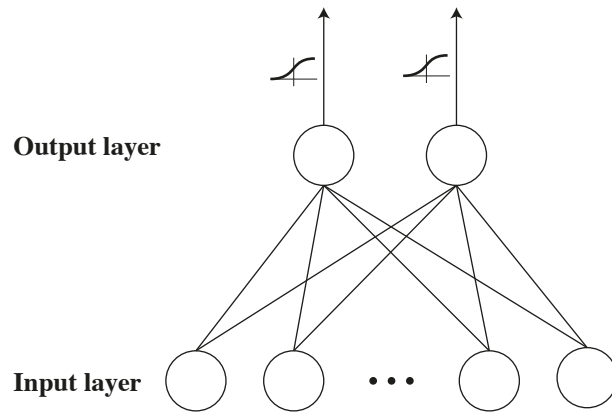


Figure 3: Schematic diagram of an adapted Rosenblatt-perceptron network

designed by Rosenblatt, it was not feasible training a model with an architecture of that kind.

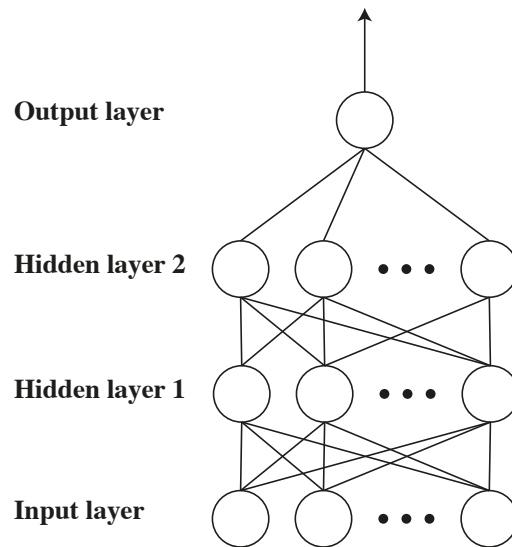


Figure 4: Schematic diagram of a multi layer perceptron (MLP)

1.2.2 Backpropagation

The procedure used to optimize the weights of a multilayer perceptron is called backpropagation. This method was first applied to neural networks by Werbos - as part of his dissertation [Werbos, 1974]- and still works the same way to this day. It is important to use a sigmoid function instead of the step function which is used in the perceptron, as this function is differentiable.

Assume a random initial distribution of the weights w_1, \dots, w_d for a network with d layers, which are the starting values for the update procedure. Let τ_i be the number of units in layer i . So layer i can be denoted as function f_{w_i} and the whole network as a chain of functions.

$$\hat{y}_i = f_{w_d}(f_{w_{d-1}}(\dots(f_{w_1}(x_i)))) = f_{w_d} \circ f_{w_{d-1}} \circ \dots \circ f_{w_1}(x_i) \quad (15)$$

The goal is to update the weights of each layer in such a way as to minimise a selected loss function $L(w_1, \dots, w_d)$. A common loss function is for example quadratic loss: $L = \frac{1}{2} \sum_{i=1}^n (\hat{y}_i - y_i)^2$. The backpropagation algorithm consists of the 3 following roughly outlined steps.

1. **forward pass:** Calculation of the result y_i for all units at the set weights
2. **backward-sweep:** Each layer's derivatives are now calculated using the chain rule, step by step - starting with the output layer. For the output layer this is:

$$\frac{\partial L}{\partial w_d} = \frac{\partial L}{\partial f_d(x_i)} \frac{\partial f_d(x_i)}{\partial w_d}$$

3. **updating:** Update the weights by an arbitrary optimization algorithm. E.g. steepest descent: $\Delta w_d = \alpha \frac{\partial L}{\partial w_d}$, where α is the learning rate

Thanks to current software, the derivatives no longer have to be calculated for each node by hand, but the standard packages are capable of "symbolic differentiation" for certain network structures, which makes training neural networks comfortable [Chollet, 2018, p. 47].

Note: While Rosenblatt for the perceptron was actually inspired by the biological structure of the neurons in the brain, more complex "networks used by engineers are only loosely based upon biology" [Hecht-Nielsen, 1988].

1.2.3 Implementation

This work uses the **Keras** R package to create neural networks, a deep learning API (Application-Programming-Interface) to deep learning backend engines, designed for R by Allaire in 2017. Keras is a so called "model-level" library, designed to set up complex neural net architectures using well arranged, high-level building blocks. As backend-engine the users are provided with TensorFlow, Theano and Microsoft Cognitive Toolkit. Using these backend engines computation may be processed seamlessly via CPU or GPU [Chollet, 2018].

2 Analysis of textbook chapters

In this first example, the chapters of individual books are classified, all of which sourced from the freely distributed Project Gutenberg. Project Gutenberg is a provider of over 60,000 free electronic books with the primary aim to “encourage the creation and distribution of eBooks”[Hart,]. The package `gutenbergr` [Robinson, 2018] preprocesses the ebooks text data and downloads it. This convenient package in combination with the free repository allows the analysis of a large number of large text documents with a secure source and a big amount of meta information. For this analysis interesting information is e.g. Gutenberg id as key variable for fast identification of books, title, author and more importantly the categorization of project Gutenberg, the so-called Gutenberg bookshelf. Table 1 lists all this information for a random sample of books downloaded from the Project Gutenberg.

Table 1: Example book corpus

gutenb. id	title	author	gutenberg bookshelf
2095	Clotelle: A Tale of the Southern States	Brown, William	African American Writers
6315	The Awakening of Helena Richie	Deland, Margaret	Bestsellers, American
6971	Judaism	Abrahams, Israel	Judaism
7635	The Disowned — Volume 05	Lytton, Edward and Baron	Historical Fiction
10319	Dave Darrin’s Third Year at Annapolis	Hancock, Harrie Irving	Children’s Book Series

In order to test NLP models, two questions can be examined on the basis of these simple text documents. Firstly, does a clustering model cluster text documents similar to a human-driven classification? In this case, this could be validated by the categorization into Gutenberg bookshelves. And secondly, how good is such a classification reproducible using a classification model?

Since very large amounts of data have to be processed and analyzed in order to model the bookshelf classification for a big collection of books, a somewhat reduced approach will be used here. One breaks down a collection of books of different categories into chapters in order to cluster respectively classify these chapters as independent documents. Similarly, you could give an example, assuming N books and separate their chapters and shuffle them, is it possible to use models to reassemble the chapters into stacks that can be assigned to individual books?

The raw text record is now transformed into the tidy format of `tidytext` (i.e. one word per row). Now it is possible to remove unnecessary stop-words. These words are predefined and can be modified for the appropriate use case if necessary. The words stem from 3 sources, “onix”, “SMART” and “snowball”, whereas the latter two are pulled from the `tm` package [Silge, 2019]. An example of stop words can be found in Table 2.

Table 2: Example of stop words from tidytext package

word	lexicon
many	SMART
mrs	onix
was	snowball
sure	onix
you're	snowball
go	onix
c	SMART
where	snowball
then	SMART

Both models use a bag of words data set as input data. This means a matrix with the documents in the M lines and the frequencies of the V used words in the entries of the columns. The dimension - i.e. the number of words in this "dictionary" V - may be reduced for two reasons. For one thing, a dimension reduction can reduce the fitting time of the model, for another thing, the diversity of the documents can be increased by skilfully reducing certain words, which occur equally frequently in all documents. Now this requires a special measure on which we decide which words to exclude. One can call this reduction of the bag of words dimensionality "embedding". In the course of this work 3 different embedding methods were tested.

1. The reduction by the words that occur with a low frequency.
2. No dimension reduction, i.e. use of the full dictionary of all occurring words.
3. The reduction with the aid of the measure *tf-idf*.

tf-idf is a combination of the term frequency and the inverse document frequency, defined as follows.

$$tf-idf(t, d) := tf(t, d) \times idf(t)$$

$$tf(t, d) := \frac{f_{t,d}}{\sum_{t_i=1}^V f_{t_i,d}}$$

$$idf(t) := \ln \left(\frac{M}{n_{d' \in t}} \right)$$

t ... term (word)

d ... document

$f_{t,d}$... frequency of term t in document d

M ... number of documents

$n_{d' \in t}$... number of documents containing term t

Even if “its theoretical foundations are considered less than firm by information theory experts”, *tf-idf* “has proved useful in text mining” [Silge and Robinson, 2017].

A Appendix

References

- [A. Gelman, 2014] A. Gelman, J. Carlin, H. S. D. D. A. V. D. R. (2014). *Bayesian Data Analysis*. Chapman and Hall/CRC.
- [AP Dempster, 1977] AP Dempster, NM Laird, D. R. (1977). Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society*.
- [Bahjat, 2006] Bahjat, F., e. a. (2006). Multivariate logistic models. *Biometrika Trust*, 93(4):1011–1017.
- [Blei, 2012] Blei, D. (2012). Probabilistic topic models. *Communications of the ACM*, 55(4):77.
- [Blei, 2003] Blei, D., N. N. J. M. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research*, 3.
- [Chollet, 2018] Chollet, F., A. W. (2018). *Deep Learning with R*. Manning.
- [Hart,] Hart, M. Project gutenber.
- [Hecht-Nielsen, 1988] Hecht-Nielsen, R. (1988). Neurocomputing: picking the human brain. *IEEE Spectrum*, 25(3):36–41.
- [Hornik, 2011] Hornik, K., B. G. (2011). topicmodels: An r package for fitting topic models. *Journal of Statistical Software*, 40(13).
- [Jordan, 1999] Jordan, M., e. a. (1999). An introduction to variational methods for graphical models. *Kluwer Academic Publishers - Machine Learning*, 37:183–233.
- [M. Steyvers, 2006] M. Steyvers, T. G. (2006). Probabilistic topic models. *Latent Semantic Analysis: A Road to Meaning*.
- [McCulloch and Pitts, 1943] McCulloch, W. and Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*.
- [Minsky and Papert, 1969] Minsky, M. and Papert, S. (1969). *Perceptrons: An Introduction to Computational Geometry*. MIT press.
- [Mukherjee, 2019] Mukherjee, S. (2019). Lecture notes for probabilistic machine learning. Duke University.
- [Powieser, 2012] Powieser, M. (2012). Latent dirichlet allocation in r. Master’s thesis, Vienna University of Economics and Business.

- [Robinson, 2018] Robinson, D. (2018). *Package 'gutenbergr'*.
- [Rosenblatt, 1958] Rosenblatt, F. (1958). The perceptron, a probabilistic model for information storage and organization in the brain. *The Psychological Review*.
- [Silge and Robinson, 2017] Silge, J. and Robinson, D. (2017). *Text Mining with R: A Tidy Approach*. O'Reilly Media.
- [Silge, 2019] Silge, J. e. a. (2019). *Package 'tidytext'*.
- [Wainwright and Jordan, 2008] Wainwright, M. J. and Jordan, M. I. (2008). Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1–2):1–305.
- [Werbos, 1974] Werbos, P. (1974). *Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences*. PhD thesis, Harvard University.

List of Figures

1	The well-established plate diagram for the standard LDA model extended by the parameter δ . The slightly bigger box represents the generative model of the corporis M documents. The smaller plate represents the iterative generation process of the N words of each document with the aid of the topics. See also "smoothed LDA model" in [Blei, 2003] for comparisons.	2
2	Schematic diagram of a simple perceptron by [Rosenblatt, 1958]	7
3	Schematic diagram of an adapted Rosenblatt-perceptron network	8
4	Schematic diagram of an multi layer perceptron (MLP) . . .	8

List of Tables

1	Example book corpus	10
2	Example of stop words from tidytext package	11