

**FOM - Hochschule für Oekonomie & Management
Hamburg**

**Master-Studiengang Big Data & Business Analytics
2. Semester**

**Development of a solution for genetic analysis of
ALL genomes by implementing Latent Dirichlet allocation**

Betreuer:

Prof. Dr. Martin Münstermann

Autor:

Jacqueline Franßen

Matrikel-Nr: 496804

2. Fachsemester

Hamburg, den 25.07.2019

Contents

1	Abstract	1
2	Introduction	2
3	Related work	3
4	Latent Dirichlet Allocation (LDA)	6
4.1	General description	6
4.2	Examples and possible use cases	8
4.3	Python package 'Gensim'	9
5	Acute Lymphoblastic Leukemia	10
5.1	Types of Leukemia and its causes	10
5.2	Examples for Genome Analysis: Next-Generation Sequencing (NGS) .	10
5.3	Data sources: National Center for Biotechnology Information (NCBI) and Ensembl genome browser 96	12
6	Development of a solution for genetic analysis of Acute Lymphoblastic Leukemia (ALL) genomes by implementing LDA	13
6.1	Problems and challenges of genetic analysis	13
6.2	First steps: Draft of developed solution	13
6.3	Proposed solution	13
6.4	Results	13
7	Conclusion and Outlook	14
7.1	Lessons learned	14
7.2	Conclusion	14
7.3	Outlook	14

List of Figures

7.1	xxx	14
-----	---------------	----

List of Tables

1 Abstract

1

2 Introduction

3 Related work

Zhao et al. ¹ describe how topic modeling can be used to analyze NGS. By implementing topic modelling, text corpus are generated ².

In the beginning of every genome analysis, there are several important questions to ask. Jurca et al. ³ recommend to ask the following questions: What are the top studied genes in breast cancer? How regulated blood cancer research is in each country? Which countries have studied the largest number of breast cancer? Which are the popular genes mentioned together by countries every year? Where do key genes lie in the soft clusters?

Jurca et al. describe a process to use large-scale text analysis of biomedical abstracts in order to generate new hypothesis about cancer biomarkers ⁴. The target is to develop a data mining methodology that patterns in genes associated with cancer. By analyzing disease-specific gene expression data, experimental data is being checked whether a gene has indeed been upregulated or downregulated with respect to a disease.

According to Xu et al. ⁵, micro Ribonucleic Acid (MIRNA)s build a class of 17-27 nucleotides single-stranded Ribonucleic Acid (RNA) molecules that regulate gene expression post-transcriptionally. In the described text-mining process, Xu et al. identified nine MIRNAs in bladder cancer and adopted protein-protein interaction sites between these miRNAs and target genes. The results of the analyzation process lead to two relationship types between bladder cancer and its MIRNA: casual and unspecified.

Topic modelling is not only used to analyze relationships between genomes but also to

¹[Zhao et al. 2016]

²[Zhao et al. 2016]

³[Jurca et al. 2016]

⁴[Jurca et al. 2016]

⁵[Xu et al. 2013]

improve diagnoses for stroke disease. Djatna et al.⁶ describe an 'Intuitionistic Fuzzy Based Decision Tree' to diagnose different types of stroke disease. To be precise, the different types of stroke diseases can be calculated by Hamming distance. The term 'Fuzzy logic' means logic that underlies the reasoning of data by way of precise estimates. It is the fastest way to map input space into output space using a degree of membership.

Lloret et al.⁷ built an automatic summarization algorithm for literature. It can include three steps: First, topic identification, second topic interpretation and third summary generation. While describing the process of textual analyzation, Lloret et al. mention a specific term: term frequency/inverse document frequency (TFIDF) which is important for topic modelling. In addition to topic-based approaches, there are graph-based approaches and discourse-based approaches. Graph-based approaches implicate nodes that represent text elements and the edges/links refer to synonymy⁸. Discourse-based approaches include Rhethorical Structure Theory (RST), Hidden-Markov-Models (HMM), RST or Bayesian models (BM).

Yang et al.⁹ describe a process of 'constructing a database for relations between human copy number variant (CNV)s and human genetic disease via systematic text mining'. In general, CNV can cause disease by gene dosage, disruption, fusion or other genetic position effects. To be more precise, there can CNV can lead to two types of autosomal variants: They can either cause deletion or amplification of the long or broken arm region of chromosomes 1-22 or can build multiples of chromosomes 1-21 (e.g. as in disease trisomy 21).

According to their article, Yang et al. used CNV database which linked the CNV information to the NCBI Gene and Ontology database. In their article, Yang et al.¹⁰ mention three steps in the text mining process. First, during the pre-processing step, unstructured fields are split into separated sentences by using Natural Language Toolkit (NLTK), a python package¹¹. After that, in the named entity recognition (NER) step, all disease mentions within DNorm system, such as MeSH IDs are recognized. In the

⁶[Djatna, Hardhienata, and Masruriyah 2018]

⁷[Lloret and Palomar 2012]

⁸[Lloret and Palomar 2012]

⁹[Yang et al. 2018]

¹⁰[Yang et al. 2018]

¹¹[Natural Language Toolkit — NLTK 3.4.1 documentation 2019]

third step, Relation extraction (RE), the positions in sentences and entities are compared to generate instances that consist of two candidate entities within one single sentence.

Yang et al. mention two more processing methods: Parallel Processing and Post Processing which includes data cleaning and statistics ¹². The term 'data cleaning' is explained as 'de-duplicating data after each step of the process to reduce repetitive operations and prevent statistical errors'. This is a very useful step in biomedicine since as a common problem, biomedical databases contain errors. For that reason, users can give feedback through a feedback mechanism to improve the quality of the databases.

¹²[Yang et al. 2018]

4 LDA

4.1 General description

According to Jurca et al. ¹ the text mining process can be divided into four steps: First, the information has to be retrieved by user queries (Information Retrieval (IR)). Second, different vocabularies and ontologies have to be integrated (NER). Third, during Information extraction (IE), relationships between biological entities in the texts are extracted by either use co-occurrence processing or Natural Language Processing (NLP). Last, there has to be gained biologically meaningful knowledge about how biological entities are related by using Knowledge Discovery (KD) methods.

Moreover, there can be distinguished between three types of clustering: hard clustering, hierarchical clustering and soft clustering. Hard clustering describes the process of separating items into distinct groups where each item is exactly in one cluster. Hierarchical clustering implicates single-link (how similar the items are to one another) and complete-link (how dissimilar the items are). Soft clustering means that items cannot be distinctly separated into clusters and partly are member of two or more clusters at a time ².

Besides, Djatna et al. ³ mention data mining techniques, such as Classification and Regression Tree (CART), Iterative Dichotomized 3 (ID3), Decision Tree (DT) and two classification techniques: Principal Component Analysis (PCA) and LDA.

LDA was developed by David Blei et. al in the year 2003 and is a clustering algorithm for text mining. It counts to the most popular topic modelling algorithms ⁴. According to ⁵, topic modelling requires of a number documents which represent each of them

¹[Jurca et al. 2016]

²[Jurca et al. 2016]

³[Djatna, Hardhienata, and Masruriyah 2018]

⁴[Zhao et al. 2016]

⁵[Zhao et al. 2016]

a mixture of latent topics. Moreover, each topic is expressed by a distribution of words. During LDA, two relationships are analyzed: First, the relationship between documents and words, also called 'per-document topic distributions'. Second, the relationship between words and topics ('per-topic word distributions'). To measure the relationships exactly and to make inference about topics and documents for text mining, probability matrices are calculated.

LDA is a probabilistic (Bayesian) model of text documents ⁶. The idea of LDA is to define a document as a collection of k topics. Each topic defines a multinomial distribution over a vocabulary which is drawn from a dirichlet.

Park et al. ⁷ define topic models as follows: Documents are no longer a collection of words, but a collection of topics. Furthermore, LDA is a generative topic model which uses a dirichlet parameter (also called dirichlet prior) to model documents. By changing the dirichlet prior, the number of topics that the model assigns to each word and document can be controlled. To be more precisely, a small dirichlet prior means a small number of topics assigned to each word. By increasing the dirichlet prior, the distribution of topics to each word rises. Moreover, the dirichlet parameter is obtained for each document and can be fitted using a maximum or estimated likelihood. If the dirichlet parameter does not fit, the gain in computational efficiency is obtained. Otherwise, there is no advantage (when dirichlet parameter fits well).

What is more, every term has a probabilistic relationship to every document. The topic model probabilities are stored as term relationships in thesaurus. The term frequencies are stored in the document index.

Process of LDA The process of LDA can be briefly described with the following steps ⁸: First, there are M documents in the corpus. Second, each document j has N_j words. In the next step, the observed value w_{ji} describes the appearance of a word i in a document j . In addition to that, all words will be clustered into K topics which are defined as object classes. Finally, each topic k is modelled as a multinomial distribution over the codebook.

⁶[Hoffman, Bach, and Blei 2010]

⁷[Park and Ramamohanarao 2009]

⁸[Wang and Grimson 2008]

As reported by Wang et al.⁹, LDA is a language model which clusters co-occurring words into topics. Moreover, documents are described as 'bag of words'. In order to analyze documents, users need to define the meaning of documents in vision problems. Wang et al.¹⁰ describe a special form of LDA: Spatial LDA. It encodes spatial structure among visual words, assuming the partition of words into documents is known a priori.

Extension of LDA are author-topic model, dynamic-topic model and correlated topic model. Wang et al. refer to the dynamic-topic model while describing how visual words are clustered into topics which correspond to object classes.

Sensitivity of LDA for information retrieval

4.2 Examples and possible use cases

Zhao et al. describe the process of analyzing genomes as follows: First, each document corresponds to one of the total number of Desoxyribonucleic acid (DNA) strains. Second, all documents had the same number of words. Third, the distribution of words for topics as well as the distribution of topics in documents were described by random variables obeying Dirichlet distributions with parameters α and β . After that, nucleotides and their orders in NGS sequences could be treated as words and the genetic information in sequences was translated and exhibited as a 'bag of words'¹¹. By using the strain-topic matrix derived from topic modelling, relationships or similarities between the strains serotypes can be found out.

Hoffmann et al.¹² describe the development of an online variational Bayes algorithm for LDA which is based on stochastic optimization with a natural gradient step. This step converges to a local optimum of the variational Bayes objective function. To be more precise, Bayesian models provide a natural way to encode assumptions about observed data. There can be distinguished between two approaches: First, sampling approaches are based on Markov Chain Monte Carlo (MCMC) sampling. Here,

⁹[Wang and Grimson 2008]

¹⁰[Wang and Grimson 2008]

¹¹[Zhao et al. 2016]

¹²[Hoffman, Bach, and Blei 2010]

a Markov chain defines the stationary distribution. Second, there are optimization approaches which are usually based on variational inference MCMC. In this case, variational Bayes optimizes the simplified parametric distribution.

images and schemes

4.3 Python package 'Gensim'

5 Acute Lymphoblastic Leukemia

5.1 Types of Leukemia and its causes

According to Jurca et al. ¹, cancer is the result of damage, especially of mutations to cell's DNA which leads to a cell losing its normal functionality and gains the ability to indefinitely multiply until normal tissue functions are impaired. This is also why malignant cancer is distributing so fast. Besides, each patient develops a different set of cancerous mutations in various genes which lead to multiple subtypes of cancer. Furthermore, some genes can be up-regulated (which means that they are transcribed more and are expressed), down-regulated (which means that they are not expressed) or can be co-expressed (which means that they are expressed at the same time ²).

As stated by Montano et al. ³, ALL is a malignant disorder originating from hematopoietic B-/T-cell precursors which are characterized by marked heterogeneity at molecular and clinical levels. There are many approaches to analyze these precursors, such as analyzing targeting of transcriptional factors (PAX5) which are involved in the pathogenesis of B-ALL. Other therapeutic and clinical approaches are genome editing techniques, i.e. the design of new therapies (Chimeric Antigen Receptors (CAR)s) and the study of genes involved in the evolution of pathogenesis.

5.2 Examples for Genome Analysis: NGS

NGS refers to post-Sanger sequencing methods ⁴. Since NGS produces large volumes of sequence data it might be very useful to use topic modelling techniques

¹[Jurca et al. 2016]

²[Jurca et al. 2016]

³[Montaño et al. 2018]

⁴[Zhao et al. 2016]

to maintain the flexibility for the level of resolution required for given experiments. According to Gasperskaja et al. ⁵, NGS does not require a priori knowledge about genomic feature, it only requires a low amount of DNA or RNA as input.

The step before analyzing two or more (multiple) genomes is called alignment which includes a comparison of two genomes. There are many different types of alignments, but Zhao et al. refer to the Multiple Sequence Alignment (MSA) by describing Multiple Sequence Comparison by Log- Expectation (MUSCLE) and CLUSTAL.

Gasperskaja et al. ⁶ mention an important question which should be asked before every genome analysis: 'Is the variance pathogenic?' and whether there is any relationship between genotype and phenotype which means that it can lead to disease or can cause a number of disorders. Moreover, there can be distinguished between beneficial (Single Nucleotide Polymorphism (SNP)) and pathogenic (nonsense variant) single nucleotide changes, large microscopically visible or chromosomal aberration. To find out whether a genome mutation is pathogenic, Gasperskaja et al. ⁷ explain that substantial information about functional genomics can be found through analysis of messenger RNA (MRNA) or complementary Ribonucleic Acid (CRNA) (which is a copy from MRNA by reverse transcription Polymerase Chain Reaction (PCR)). Methods to measure RNA expression are for example: Serial Analysis of Gene Expression (SAGE) or Quantitative real-time Polymerase Chain Reaction (QPCR). By using complementary Desoxyribonucleic acid (CDNA) microarray assays important genome-wide information about changes of gene expression in various cell lines can be found out.

As claimed by Montano et al. ⁸, the development of NGS techniques implicates vast amount of data which need to be translated. One important question is to find out how the genotype (the genetic expression of a biological attribute) influences the phenotype. By integrating genome editing systems into their research process, investigators are able to manipulate virtually any gene in a diverse range of cell types and organisms. To give an example, Gasperskaja et al. ⁹ and Montano et al. ¹⁰ describe

⁵[Gasperskaja and Kučinskas 2017]

⁶[Gasperskaja and Kučinskas 2017]

⁷[Gasperskaja and Kučinskas 2017]

⁸[Montaño et al. 2018]

⁹[Gasperskaja and Kučinskas 2017]

¹⁰[Montaño et al. 2018]

Clustered Regularly Interspaced Short Palindromic Repeats Cas-9 (CRISPR-CAS9). In fact, this genome analysis includes generating a direct cut in the double strand of DNA by Cas9 nuclease. Cas9 is driven by a single 20-nucleotide RNA strand which marks the direct breakpoint. After cutting the DNA, the repair machinery of the host cell repairs errors and promotes a modification of the original sequence by a mutation (e.g. insertion, deletion, inversion). But why should be used such a complicate process or why should we change the shape of DNA? In the opinion of Montano et al.¹¹, the use of genetically modified cell lines and animal models help us to better understand the functions of genes and their pathogenesis in diseases, such as cancer.

5.3 Data sources: NCBI and Ensembl genome browser

96

¹¹[Montaño et al. 2018]

6 Development of a solution for genetic analysis of ALL genomes by implementing LDA

6.1 Problems and challenges of genetic analysis

Quality of data Is the data complete, which includes that it contains all required genomes which can cause ALL.

our assumptions can lead to false content or solutions

6.2 First steps: Draft of developed solution

To get useful data, the NCBI ¹ was used to get all currently detected mutations of genomes which may cause LDA.

The first idea was to build a parsing application, which iterates over the found 582 genomes. After the iteration, it compares the oncogenes with the healthy genomes and to figure out where the differences are. The results might be displayed in a diagram. It might be possible to create clusters from the differences between the two groups or practice LDA on the differences.

6.3 Proposed solution

6.4 Results

¹[Biotechnology et al. 2019]

7 Conclusion and Outlook

7.1 Lessons learned

7.2 Conclusion

7.3 Outlook

Figure 7.1: Bildunterschrift

Ein Zitat

□

Literaturverzeichnis

- Biotechnology, National Center for et al. (2019). *National Center for Biotechnology Information*. en. URL: <https://www.ncbi.nlm.nih.gov/> (visited on 05/09/2019).
- Djatna, Taufik, Medria Kusuma Dewi Hardhienata, and Anis Fitri Nur Masruriyah (Dec. 2018). "An intuitionistic fuzzy diagnosis analytics for stroke disease". en. In: *Journal of Big Data* 5.1, p. 35. ISSN: 2196-1115. DOI: 10.1186/s40537-018-0142-7. URL: <https://journalofbigdata.springeropen.com/articles/10.1186/s40537-018-0142-7> (visited on 04/07/2019).
- Gasperskaja, Evelina and Vaidutis Kučinskas (2017). "The most common technologies and tools for functional genome analysis". In: *Acta Medica Lituanica* 24.1, pp. 1–11. ISSN: 1392-0138. DOI: 10.6001/actamedica.v24i1.3457. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5467957/> (visited on 05/09/2019).
- Hoffman, Matthew, Francis R. Bach, and David M. Blei (2010). "Online Learning for Latent Dirichlet Allocation". In: *Advances in Neural Information Processing Systems* 23. Ed. by J. D. Lafferty et al. Curran Associates, Inc., pp. 856–864. URL: <http://papers.nips.cc/paper/3902-online-learning-for-latent-dirichlet-allocation.pdf> (visited on 04/19/2019).
- Jurca, Gabriela et al. (Dec. 2016). "Integrating text mining, data mining, and network analysis for identifying genetic breast cancer trends". en. In: *BMC Research Notes* 9.1. ISSN: 1756-0500. DOI: 10.1186/s13104-016-2023-5. URL: <http://bmcresearchnotes.biomedcentral.com/articles/10.1186/s13104-016-2023-5> (visited on 04/06/2019).
- Lloret, Elena and Manuel Palomar (Jan. 2012). "Text summarisation in progress: a literature review". In: *Artificial Intelligence Review* 37.1, pp. 1–41. ISSN: 1573-7462. DOI: 10.1007/s10462-011-9216-z. URL: <https://doi.org/10.1007/s10462-011-9216-z>.

- Montaño, Adrián et al. (July 2018). "Targeted genome editing in acute lymphoblastic leukemia: a review". In: *BMC Biotechnology* 18.1, p. 45. ISSN: 1472-6750. DOI: 10.1186/s12896-018-0455-9. URL: <https://doi.org/10.1186/s12896-018-0455-9> (visited on 04/19/2019).
- Natural Language Toolkit — NLTK 3.4.1 documentation (2019). URL: <https://www.nltk.org/> (visited on 05/09/2019).
- Park, Laurence A. F. and Kotagiri Ramamohanarao (2009). "The Sensitivity of Latent Dirichlet Allocation for Information Retrieval". en. In: *Machine Learning and Knowledge Discovery in Databases*. Ed. by Wray Buntine et al. Lecture Notes in Computer Science. Springer Berlin Heidelberg, pp. 176–188. ISBN: 978-3-642-04174-7.
- Wang, Xiaogang and Eric Grimson (2008). "Spatial Latent Dirichlet Allocation". In: *Advances in Neural Information Processing Systems 20*. Ed. by J. C. Platt et al. Curran Associates, Inc., pp. 1577–1584. URL: <http://papers.nips.cc/paper/3278-spatial-latent-dirichlet-allocation.pdf> (visited on 04/19/2019).
- Xu, Yahong et al. (Sept. 2013). "RETRACTED ARTICLE: Integrated gene network analysis and text mining revealing PIK3R1 regulated by miR-127 in human bladder cancer". In: *European Journal of Medical Research* 18.1, p. 29. ISSN: 2047-783X. DOI: 10.1186/2047-783X-18-29. URL: <https://doi.org/10.1186/2047-783X-18-29>.
- Yang, Xi et al. (Dec. 2018). "Constructing a database for the relations between CNV and human genetic diseases via systematic text mining". en. In: *BMC Bioinformatics* 19.S19, p. 528. ISSN: 1471-2105. DOI: 10.1186/s12859-018-2526-2. URL: <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-018-2526-2> (visited on 04/11/2019).
- Zhao, Weizhong et al. (May 2016). "A novel procedure on next generation sequencing data analysis using text mining algorithm". In: *BMC Bioinformatics* 17.1, p. 213. ISSN: 1471-2105. DOI: 10.1186/s12859-016-1075-9. URL: <https://doi.org/10.1186/s12859-016-1075-9>.

Rechtsquellenverzeichnis

Bundesdatenschutzgesetz, BDSG, 1990, zuletzt geändert 2009

Gesetz gegen unlauteren Wettbewerb, UWG, 2004, zuletzt geändert 2013

Richtlinie 95/46/EG des Europäischen Parlaments und des Rates vom 24. Oktober 1995 Nr. L 281/31 zum Schutz natürlicher Personen bei der Verarbeitung personenbezogener Daten und zum freien Datenverkehr

Telemediengesetz, TDG, 2007, zuletzt geändert 2010

Ehrenwörtliche Erklärung

Hiermit versichere ich, dass die vorliegende Arbeit von mir selbstständig und ohne unerlaubte Hilfe angefertigt worden ist, insbesondere dass ich alle Stellen, die wörtlich oder annähernd wörtlich aus Veröffentlichungen entnommen sind, durch Zitate als solche gekennzeichnet habe. Ich versichere auch, dass die von mir eingereichte schriftliche Version mit der digitalen Version übereinstimmt. Weiterhin erkläre ich, dass die Arbeit in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde / Prüfungsstelle vorgelegen hat. Ich erkläre mich damit nicht einverstanden, dass die Arbeit der Öffentlichkeit zugänglich gemacht wird. Ich erkläre mich damit einverstanden, dass die Digitalversion dieser Arbeit zwecks Plagiatsprüfung auf die Server externer Anbieter hochgeladen werden darf. Die Plagiatsprüfung stellt keine Zurverfügungstellung für die Öffentlichkeit dar.

Ort, Datum (Vorname Nachname)