

Generative probabilistic model with Dirichlet prior distribution for similarity analysis of research topic

John Milyahilu[†], Jong Nam Kim^{††}

ABSTRACT

We propose a generative probabilistic model with Dirichlet prior distribution for topic modeling and text similarity analysis. It assigns a topic and calculates text correlation between documents within a corpus. It also provides posterior probabilities that are assigned to each topic of a document based on the prior distribution in the corpus. We then present a Gibbs sampling algorithm for inference about the posterior distribution and compute text correlation among 50 abstracts from the papers published by IEEE. We also conduct a supervised learning to set a benchmark that justifies the performance of the LDA (Latent Dirichlet Allocation). The experiments show that the accuracy for topic assignment to a certain document is 76% for LDA. The results for supervised learning show the accuracy of 61%, the precision of 93% and the f1-score of 96%. A discussion for experimental results indicates a thorough justification based on probabilities, distributions, evaluation metrics and correlation coefficients with respect to topic assignment.

Key words: Gibbs Sampling, Corpus, Machine Learning, Probabilistic Model, Topics

1. INTRODUCTION

The standard way for document searching on internet is through keywords, phrases and clauses. Search engines like Bing, Yahoo, Naver, Google and others of alike routinely use machine learning algorithms to provide information requested by the user at an outstanding performance. The tools mentioned above perform better but when the process of searching involves large documents, they may get poor performance [1]. On the other hand, topic modeling and classification help a person to identify the topics and main themes for large collection of documents.

Topic modeling is a technique that uses stat-

istical models to discover topics and themes that occur in a collection of documents known as corpus [2]. Topic models are also known as probabilistic models that use statistical algorithms for discovering the latent semantic structures of a corpus [3]. They deal with categorizing texts into organized groups for the purpose of getting insights about them, discovering main topics, and describing themes [4]. Topic modeling with LDA gives a summary about the theme of a corpus. It works under an assumption that each document has a mixture of topics and each topic has a collection of words [5]. The applicable techniques for topic modeling include LDA based on probabilistic graphical models while Latent Semantic Analysis (LSA), Latent

* Corresponding Author : Jong Nam Kim, Address : (48513) Yongso-ro 45, Nam-gu, Busan, Korea, TEL : +82-51-629-6259, FAX : +82-51-629-6263, E-mail : jnkim1225@gmail.com

Receipt date : Apr. 7, 2020, Revision date : Apr. 14, 2020
Approval date : Apr. 16, 2020

[†] Dept. of IT Convergence & Application Eng. Pukyong National University
(E-mail : johnmilyahilu@gmail.com)

^{††} Dept. of IT Convergence & Application Eng. Pukyong National University

* This research was supported by HK+ Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education, Science and Technology (NRF - 2017S1A6A3A01079869)

Semantic Indexing (LSI) and Non-negative Matrix Factorization (NMF) base on linear algebra [6].

The motivation of this work bases on the challenge that mega projects are subjected to publications of journal papers via reputable ones. This model has the capacity to assign a particular topic to a document within a corpus based on the objectives of the project and calculate the correlation coefficients among the documents. The model serves as a proof to some projects whether they correlate with the papers published through different journals with respect to the project's goals and objectives. Basically, the proposed work will be able to accept or reject some papers that belongs or do not belong to a certain project's objectives. In this work we propose a generative topic model that is capable of jointly modeling the words in a corpus arising from a mixture of topics within a corpus. This model generates both a per-document topic distribution and a per-topic word distribution that is used to determine the posterior distribution. Using posterior probabilities assigned to documents in a corpus, we determine the relationship among documents in the corpus and determine whether they belong to a certain project or not. The reason of the better performance from the proposed work is the highest probability assigned to the corresponding word in a document term matrix.

This work is organized into five sections. The second section presents the related works and the third section describes the proposed work. Section four gives experimental results and discussion. Lastly, section five presents concluding remarks.

2. RELATED WORKS

Topic modeling and classification are not new fields, but their significant contribution to the study of natural language processing is pivotal up to recently. There are several approaches that are used to model topics such as dimensionality reduction techniques, unsupervised learning such as

clustering and form of tagging techniques [7]. Due to their essence, research works for topic modeling and classification are tremendously increasing [8]. Notable progress has been made on solving topic modeling problems by using the popular TF-IDF (Term Frequency-Inverse Document Frequency) algorithm that employs words or terms and their respective frequencies for each document in a corpus [9]. This popular algorithm has some shortcomings such as small reduction in descriptive length for documents and limited discovery of topic and sub-topic in a collection [10].

In addressing the shortcomings of the previous schema, Deerwester et al. introduced the most eminent schema known as LSI [11]. It uses singular value decomposition as a dimension reduction algorithm which creates a document term matrix that captures most of the variance within a corpus. The schema was criticized by many scholars while substantiating the weaknesses of the schema that were outlined by Hofmann [12]. The model could handle multiple documents however, it had some shortcomings which are overfitting caused by the increase in number of parameters due to the size of document and assigning probability to a document outside the training set was not clearly put into description. Blei et al. introduced LDA algorithm which performs better than other described algorithms on topic modeling because it employs variational methods and expectation maximization algorithms for estimating empirical Bayes parameters [10].

LDA model performs better than TF-IDF schema due the use of variational Bayesian estimation that overcomes the shortcomings mentioned in the previous paragraphs of this section. LDA model as unsupervised learning technique is a basis for text classification and performs better than supervised learning techniques in terms of accuracy. There are several classifiers which are used in text analysis such as Naïve Bayes classifier, Support Vector Machines (SVM), Decision Trees, and Nearest

neighborhood. Among others, Naïve Bayes classifier has better performance than SVM in text categorization as described in [13]. Currently, text classification is efficiently done by deep learning algorithms including convolutional neural networks as described in [14]-[17].

We are interested in determining a probabilistic model for topic modeling that assigns probability to each document in a corpus and has an ability to determine the relationship among documents by assigning similarity coefficients.

3. PROPOSED WORK

LDA is a supervised machine learning algorithm produces two matrices which are a document-topic matrix and a topic-terms matrix. We certainly run unsupervised method and transform the data using the algorithm for better descriptions of the results. We subsequently examine how a method performs on the transformed data, and compare the results based on precision, f1-score and recall when used on the same corpus transformed by LDA and TF-IDF algorithms [10],[16]. The LDA posterior cannot determine document similarity but might be used to determine topic similarity. The assigned probabilities can therefore be used to determine the document similarity. The proposed work helps to assign topic on different documents in a corpus by using the posterior probabilities that are obtained from the model. LDA posterior probabilities are obtained from the multinomial distribution using conjugate Dirichlet prior. We use the Dirichlet prior because our problem is a multivariate and is difficult to do direct sampling but its probability distribution belongs to the same family of the posterior distribution.

The procedure of the proposed algorithm is summarized in Fig. 1 whereby the first block with Raw data represents a corpus made up by some texts. The second block with the Preprocessing removes spaces, numbers, special characters and stop words in the document. The third block with

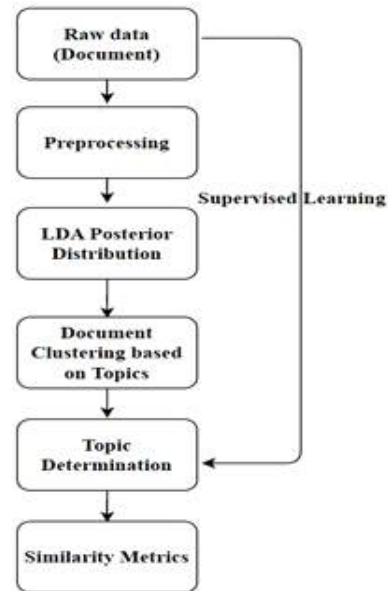


Fig. 1. Procedure of the proposed algorithm.

LDA Posterior Distribution calculates posterior distribution from text of the Preprocessing block. It contains Gibbs sampler which is the conjugate prior of the posterior distribution (multinomial). The conjugate prior is the exact distribution of the posterior distribution. The fourth block with Document Clustering designates topics randomly after obtaining the posterior probabilities from LDA Posterior Distribution block. The fifth block with Topic Determination assigns topics that can be done by both supervised learning and LDA algorithms. The last block with Similarity Metrics determines correlation coefficients from the clustered documents in the Document Clustering block.

The LDA Posterior block receives output from the Preprocessing block and uses Gibbs sampler to generate the outputs that are finally used in Document Clustering block. Gibbs sampling is a Markov Chain Monte Carlo (MCMC) algorithm for obtaining sequence of observations which are approximately from the multivariate distribution. We propose Gibbs sampling technique because direct sampling is difficult to perform i.e. the probability density functions of the samples are not analytical.

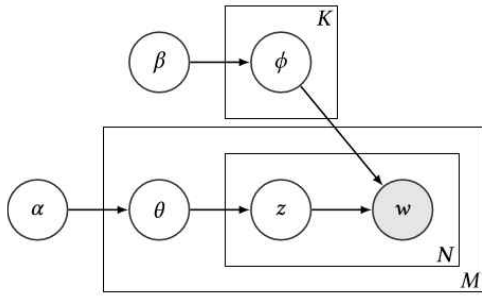


Fig. 2. Graphical representation of LDA.

From Fig. 2 we determine the posterior distribution by letting the parameters for the Gibbs sampler as follows; prior parameter alpha (α) be the per-document topic proportions, z be the per-word topic proportion, w be the observed word in a document, T be per-topic word proportions, beta (β) be topic hyperparameter, and N_d be a number of documents in a corpus.

That means, there are T topics $\phi_{1,\dots,T}$ that are shared among all documents, and each document D_j in a corpus D . The corpus is assumed as a mixture of topics indicated by θ_j . We can first sample a topic assignment $z_{j,t}$ from the topic proportions to generate the words for each document. We assume that the generative model has T topics, a corpus D of N_d documents and a vocabulary consisting of V unique words. We choose the Dirichlet parameter β as the prior for distribution ϕ_i of the generative model and present the process using the following pseudocode supported by the graphical model representation.

Pseudocode of Gibbs Sampling algorithm

- For $i \in [1, \dots, T]$

Generate the multinomial distribution over topics in a corpus $\phi_i \sim \text{Dirichlet}(\beta)$

- For $j \in [1, \dots, M]$,

Generate multinomial distributions over the documents in a corpus $\theta_j \sim \text{Dirichlet}(\alpha)$

For $t \in [1, \dots, N_d]$

Generate topic $z_{j,t} \sim \text{Multinomial}(\phi_j)$

Generate term $w_{j,t} \sim \text{Multinomial}(\phi_j)$

Calculate the correlation coefficients for each document $\theta_{i,j}$ in a corpus D .

The graphical representation of the LDA algorithm is given in Fig. 2. The outer block represents documents, while the inner block represents the repetitive choices of topics and words within a document.

The LDA algorithm in Fig. 2 is divided into three parts which are unigram model, mixture of unigrams and probabilistic latent semantic indexing model. In the inner block of this figure which constitutes of w and N from a corpus M , the terms(w) of every document are drawn independently from a multinomial distribution expressed by Eq. (1). This stage represents the unigram model within the LDA model for discrete text data.

$$p(w) = \prod_{j=1}^{N_d} p(w_n), \tag{1}$$

The second stage is the combination of the previous stage and z that makes a mixture of unigrams stage. In this stage each document is generated by choosing a topic z and then generate N terms from multinomial distribution $p(w|z)$ as shown in Eq. 2.

$$p(w) = \sum_z p(z) \prod_{n=1}^N p(w_n|z), \tag{2}$$

The third stage is the probabilistic latent indexing, which is the combination of the first stage, the second stage with the document D_j from θ . In this stage, the mixture of topics in θ and the term w_n are conditionally independent for unknown topic z as expressed in Eq. 3.

$$p(\theta, w_n) = p(\theta) \sum_z p(w_n|z) p(z|\theta), \tag{3}$$

where θ is a document chosen in a multinomial distribution with a probability density function of ϕ_j . This stage determines the possibility that a document contains multiple topics because of $p(z|\theta)$ acts as a mixture of weights of topics in a

document. The probability distribution of θ is a multivariate Dirichlet distribution given by Eq. 4.

$$p(\theta|\alpha) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \theta_1^{\alpha_1-1} \dots \theta_k^{\alpha_k-1}, \tag{4}$$

where $\Gamma(\cdot)$ is the Gamma function and α is the parameter within $(k-1)$ simplex or k vector also known as a multinomial distribution of k words.

Lastly, we use Pearson correlation coefficient in Eq. (5) to determine the relationship (similarities) among documents using the probabilities assigned to the topics of each document.

$$r = \frac{n(\sum_{i,j=1}^n D_i D_j) - \sum_i D_i \sum_j D_j}{\sqrt{n[\sum_{i=1}^n D_i^2 - (\sum_{i=1}^n D_i)^2]n[\sum_{j=1}^n D_j^2 - (\sum_{j=1}^n D_j)^2]}} \tag{5}$$

4. EXPERIMENTAL RESULTS AND DISCUSSION

The experiments of this work were conducted through a computer operating on windows 10 installed with Python and R programs that were used for data analysis. Text mining, natural language processing, machine learning and topic models' libraries were employed to analyze the performance. The dataset were the extracted abstracts from scientific papers published by IEEE journals from 2010 to 2019 in the field of "motion estimation", "image classification", "image segmentation", "object detection" and "3D reconstructions". Each field has a total 10 papers with different titles and contents while some papers share the methodologies and techniques for analysis.

The preprocessing steps involve transforming to

lower case, removing special symbols, numbers, punctuations, white spaces and general errors via tm library which is related to text mining. The initial entry argument in the topic model function is the document-term matrix associated with parameters such as number of clusters. We set a few initial steps of a random walk (Gibbs Sampler) and skip the steps that are not related to the prior distribution through the burn-in period and perform 2000 iterations dividing them into four steps in order to overcome multicollinearity. The random walk finds only a local optimal solution that is obtained by implementing different settings on the parameters for many times.

The algorithm with a Gibbs sampler performs better than the supervised learning approach whose results are displayed in different tables showing that this method is robust and efficient on topic modeling problems. Based on the document term matrix, the assigned topics are image segmentation, image classification, 3D reconstruction, motion estimation and object detection. The following tables show the most frequently used terms in each topic, document-topic distribution, probabilities, evaluation metrics and correlation coefficients.

Table 1 describes the topic word distribution with the most frequently used terms for each assigned topic to every document in a corpus. In the first topic the most frequent words are "image", "methods", "use" and "segmentation." If we assume that the word "use" and "method" are stop words then, the only significant words could represent the topic of "image segmentation". The same procedure applies to the rest of the topics starting from the second topic to fifth one that ignores the

Table 1. Showing topic-word distribution for each assigned topic

Term	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5
1	image	network	paper	algorithm	learn
2	method	feature	reconstruction	propose	detect
3	use	model	present	search	object
4	segment	classify	result	motion	deep

Table 2. Document–topic distribution

Document	Topic	Document	Topic
abstract1.txt	2	abstract30.txt	2
abstract10.txt	3	abstract31.txt	4
abstract11.txt	2	abstract32.txt	4
abstract12.txt	2	abstract33.txt	4
abstract13.txt	1	abstract34.txt	4

Table 3. Topic probabilities by documents (abstracts)

Abstracts	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5
20	0.44	0.13	0.11	0.11	0.21
21	0.51	0.08	0.15	0.08	0.18
25	0.09	0.08	0.11	0.60	0.11
27	0.11	0.13	0.12	0.59	0.07
28	0.12	0.14	0.12	0.55	0.06
39	0.11	0.16	0.09	0.08	0.56
42	0.12	0.17	0.06	0.09	0.56

term “learn” and “deep” to represent the topic “object detection”.

Results in Table 2 show that, from documents named abstract10–13.txt, there are only two abstracts that are correctly assigned to a corresponding topic and the rest were wrongly assigned, abstracts30–34.txt were correctly assigned except for the document named abstract30.txt. The topic assignment accuracy is equivalent to 76% for 50 documents involved in the experiment.

The highest probabilities are observed for docu-

ments assigned to topic 4 and 5 while topic 1 and 3 contained smallest probability values. The model performed poorly on assigning topics to documents named as abstracts1–10.txt whereby the correctly assigned documents were only four.

Table 4 displays the evaluation metrics for a supervised learning approach classifying texts per predefined labels or topics. We randomly divided the training dataset into two parts. The one is 25% of the data for training the classifier and the other is 75% of the dataset for testing it [13]. Using Naïve Bayes classifier, we obtained 61% of accuracy, which is correctly assigned to a topic. Additionally, we got 93% of precision that the classifier correctly tagged each document to a topic. The fraction of the relevant documents that are successfully retrieved is 100% and f1-score is 96%.

Table 5 shows the similarity between topics assigned to different documents in a corpus. Documents assigned to “image classification” and “image segmentation” were closely related to each other while documents assigned to “image classification” and “3D reconstructions” were not related to each other. The topics assigned as “motion estimation” and “object detection” indicate positive correlation coefficients while the documents assigned to “3D reconstructions” are negatively correlated to the documents assigned to “Image segmentation”. From the experimental results above, it is noted that the proposed algorithm is efficient

Table 4. Evaluation metrics for text classification

Topics	Precision	Recall	F1-score		Support
3D Reconstruction	0.69	0.64	0.67		14
Convolutional Neural Networks	0.67	0.67	0.67		6
Classification	0.42	0.42	0.42		12
General Statements	0.30	0.41	0.35		17
Motion Estimation	0.87	0.93	0.90		14
Object Detection	0.60	0.43	0.50		14
Segmentation	0.62	0.45	0.53		11
Overall Evaluation Metrics	Precision	Recall	F1-Score	Accuracy	13
	0.93	1.00	0.96	0.61	101

Table 5. Document similarity based on Pearson correlation coefficient

	A1	A5	A11	A15	A21	A25	A31	A35	A41	A45
A1	1.00	0.71	0.78	-0.42	-0.11	-0.02	0.16	0.32	0.48	0.54
A5	0.71	1.00	0.75	-0.82	0.22	-0.81	-0.84	-0.71	0.60	-0.02
A11	0.78	0.75	1.00	0.99	-0.48	-0.47	-0.83	-0.51	0.02	-0.08
A15	-0.42	-0.82	0.99	1.00	0.71	0.68	0.39	0.24	-0.52	0.59
A21	-0.11	0.22	-0.48	0.71	1.00	0.80	-0.03	-0.42	0.51	0.48
A25	-0.02	-0.81	-0.47	0.68	0.80	1.00	-0.56	0.22	0.45	0.66
A31	0.16	-0.84	-0.83	0.39	-0.03	-0.56	1.00	0.94	0.36	0.44
A35	0.32	-0.71	-0.51	0.24	-0.42	0.22	0.94	1.00	0.64	0.73
A41	0.48	0.60	0.02	-0.52	0.51	0.45	0.36	0.64	1.00	0.86
A45	0.54	-0.02	-0.08	0.59	0.48	0.66	0.44	0.73	0.86	1.00

and accurate for assigning topics to documents and obtains correlation coefficients for related and un-related documents.

5. CONCLUSION

In the paper, we proposed a generative probabilistic model with Dirichlet prior distribution for topic modeling and text similarity analysis. The characteristics of the proposed algorithm was to determine the posterior distribution based on the Gibbs sampler. The algorithm determined document similarity based on the posterior probabilities assigned to each topic of a document in a corpus. The meaning of the results is that the highest or lowest probability on a topic assigned to a document can affect the correlation coefficient among documents. The experimental results indicate that 76% of the documents in a corpus were correctly assigned to specific topics using the unsupervised learning technique. On the other hand, supervised learning technique based on a Naïve Bayes classifier achieved an accuracy of 61% which is smaller than that of the proposed algorithm. From the results, we can conclude that the proposed algorithm has good performance for assigning topics to documents and text similarity analysis.

REFERENCE

- [1] H. Yuening, J.B. Graber, B. Satinoff, and A. Smith, "Interactive Topic Modeling," *Machine Learning*, Vol. 95, pp. 423–469, 2014.
- [2] H. Jelodar, Y. Wang, C. Yuan, X. Feng, X. Jiang, Y. Li, et al., "Latent Dirichlet Allocation (LDA) and Topic Modeling: Models, Application, A Survey," *Multimedia Tools and Applications*, Vol. 78, pp. 15169–15211, 2019.
- [3] A. Beykikhoshk, O. Arandjelovic, D. Phung, and S. Venkatesh, "Discovering Topic Structures of a Temporary Evolving Document Corpus," *Knowledge and Information Systems*, Vol. 55, pp. 599–632, 2018.
- [4] K. Kowsari, K.J. Meimandi, M. Heidarysafa, S. Mendu, L. Barnes, and D. Brown, "Text Classification Algorithms: A Survey," *Information*, Vol. 10, No. 3, pp. 1–68, 2019.
- [5] N. Pladeau, and E. Davoodi, "Comparison of Latent Dirichlet Modeling and Factor Analysis for Topic Extraction," *Proceeding of the Hawaii International Conference on System Sciences*, pp. 615–623, 2018.
- [6] J. Rashid, S.M. Shah, A. Irtaza, T. Mahmood, M. Shafiq, and A. Gardezi, "Topic Modeling Technique for Text Mining Over Biomedical Text Corpora Through Hybrid Inverse Documents Frequency and Fuzzy K-Means Clustering," *IEEE Access*, Vol. 7, pp. 146070–146080, 2019.
- [7] J. Clark, and F. Provost, "Unsupervised Dimension Reduction versus Supervised Re-

- gularization for Classification from Sparse Data," *Data Mining and Knowledge Discovery*, Vol. 33, pp. 871-916, 2019.
- [8] T.R. Hannigan, R.F. Haans, K. Vakili, H. Tchalian, V.L. Glaser, M.S. Wang, et al., "Topic Modeling in Management Research," *Academy of Management Annals*, Vol. 13, No. 2, pp. 586-632, 2019.
- [9] Z. Gou, Z. Huo, K. Vakili, Y. Liu, and Y. Yang, "A Method for Constructing Supervised Topic Model Based on Term Frequency-Inverse Topic Frequency," *Symmetry*, Vol. 11, pp. 1-9, 2019.
- [10] D.M. Blei, A.Y. Ng, and M.I. Jordan, "A Latent Dirichlet Allocation," *Journal of Machine Learning Research*, Vol. 3, pp. 993-1022, 2003.
- [11] S. Deerwester, S.T. Dumais, G.W. Furnas, T.K. Landauer, and R. Harshman, "Indexing by Latent Semantic Analysis," *Journal of American Society for Information Science*, Vol. 41, No. 6, pp. 391-407, 1990.
- [12] T. Hofmann, "Probabilistic Latent Semantic Indexing," *Proceeding of International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 50-57, 1999.
- [13] M.K. Dalal, X. Chuangbai, I. Izahar, and T. Shanshan, "Automatic Text Classification: A Technical Review," *International Journal of Computer Applications*, Vol. 28, No. 2, pp. 37-40, 2011.
- [14] H. Sundus, and R. Muhammed, and M. Shaikh, "Comparing SVM and Naive Bayes Classifiers for Text Categorization with Wikitology as Knowledge Enrichment," *Proceeding of IEEE International Conference*, pp. 1-3, 2012.
- [15] M.L. Prabha, and G.U. Srikanth, "Survey of Sentiment Analysis Using Deep Learning Techniques," *Proceeding of International Conference on Innovations, in Information and Communication Technology*, pp. 1-9, 2019.
- [16] R.K. Roul, J.K. Sahoo, and K. Arora, "Modified TF-IDF Term Weighting Strategies for Text Categorization," *Proceeding of IEEE India Council International Conference*, pp. 1-6, 2017.
- [17] M. Nam, E. Lee, and J. Shin, "A Method for User Sentiment Classification using Instagram Hashtags," *Journal of Korea Multimedia Society*, Vol. 18, No. 11, pp. 1391-1399, 2015.



John Mlyahilu

2009. 12. : BS Mathematics,
University of Dares
Salaam

2014. 2. : MS Statistics, Pukyong
National University

2018. 9. ~ Now : Ph.D. Student,
Pukyong National
University

Research fields: multimedia and image processing,
computer vision, AI



Jong Nam Kim

1997. 2. : MS Information
Telecommunication,
GIST

2001. 8. : Ph.D. Mechatronics,
GIST

2001. 8. ~ 2004.2. : Researcher at
KBS

2004. 3. ~ Now : Professor at Dept. of IT Conv. & Apps
Engineering, Pukyong National
University

Research fields: video compression, image processing,
computer vision, deep learning