

UNIVERSITY COLLEGE LONDON

MASTERS THESIS

Dynamic Topic Modeling of PATSTAT Patents Using LDA

Author:
Christopher MARTIN

Supervisor:
Dr. John Shawe TAYLOR

*A thesis submitted in fulfillment of the requirements
for the degree of Masters of Science*

Machine Learning
UCL Dept. of Computer Science

August 15, 2016

Declaration of Authorship

I, Christopher MARTIN, declare that this thesis titled, “Dynamic Topic Modeling of PATSTAT Patents Using LDA” and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

UNIVERSITY COLLEGE LONDON

Abstract

Faculty Name
UCL Dept. of Computer Science

Masters of Science

Dynamic Topic Modeling of PATSTAT Patents Using LDA

by Christopher MARTIN

In this paper we evaluate the performance of a time varying family of LDA based topic models meant to capture both the underlying semantic structure of a document collection and the evolution of that structure in time. Such models are useful for illustrating changes in the use of language regarding specialized subject matter and provide a window into the progression of that change. We compare these models to traditional topic models such as LDA as a benchmark and explore the efficacy of such models in a range of applications including document classification, clustering, and influence prediction. Finally, we present results on over 18 years of patent data from the PATSTAT database across X classes of patents demonstrating interpretable trends, better document classification and clustering, and improved topic coherence.

Acknowledgements

I owe the completion of this project to the many people who have helped along the way, either directly or indirectly. To my project supervisors Christopher Grainger and Prof. John Shawe-Taylor for their continued guidance throughout the project, to my colleagues for their advice and discussion, and to my family for their moral support, I express my sincere gratitude, thank you.

Contents

Declaration of Authorship	ii
Abstract	iii
Acknowledgements	iv
1 Introduction to the Thesis Topic	1
1.1 The need for topic models	1
1.2 Primer on Latent Dirichilet allocation	1
1.3 Adding a temporal component	4
1.4 Why patents?	4
1.5 Experiments	5
Historical Topic Trend Validation	6
1.5.1 Topic Coherence	6
1.5.2 Classification	6
1.5.3 Clustering	6
2 Background Information and Theory	7
2.1 Literature Review	7
2.2 DTM model	7
2.2.1 Issues	7
2.3 DIM model	7
3 Experimental Set Up	8
3.1 Data Prep and Considerations	8
4 Experimental Results	9
4.1 DTM Results and Insights	9
4.1.1 Topics Through Time	9
validating topic histories in technology	9
4.2 DIM Results and Insights	9
4.2.1 Influence Metric	9
validating influential patents	9
correlation with forward citations	9
correlation with page-rank	9
4.3 Performance Evaluation	9
4.3.1 Classification	9
4.3.2 Clustering	9
5 Usefulness in Other Models	10
5.1 Economic Model	10
6 Conclusions and Future Work	11
6.1 Conclusions	11
6.2 Future Work	11

A Appendix Title Here	12
Bibliography	13

List of Figures

1.1	Patent114	2
1.2	Graphical model for LDA	3
1.3	wwtfTopic6	5

List of Tables

Chapter 1

Introduction to the Thesis Topic

1.1 The need for topic models

Researchers today are faced with a deluge of data. As we continue to digitize and aggregate our collective knowledge we produce ever increasing archives of information. The sheer volume and variety of forms this information may take - text, images, audio, video, social connections etc. - makes it difficult and in most cases impossible to parse manually.

This driving factor of data growth has given rise to internet giants such as Google, whose search tool helps us access and browse pre-indexed swathes of information. However in order to go beyond mere keyword searches, or link analysis, and break into the realm of understanding each document, requires a new approach to data exploration.

A powerful set of computational tools referred to as probabilistic topic models have emerged to meet this challenge. Aimed to discover and annotate large archives of documents with thematic information, topic models identify patterns that reflect the underlying topics which combined to form those documents.

Naturally, it is rare that we would know beforehand exactly what topics a given document contains, and thus topic modeling constitutes an unsupervised task. As a result, topic modeling algorithms are designed to work without prior knowledge of the topic distribution of a given document — that is, the topics are woven from the texts themselves. This makes the organization, summarization and annotation of text corpora possible at an inhuman scale. Consequently, topic models are useful in a variety of settings and have successfully been applied to web archives, news articles (Newman et al., 2006), and academic literature (Steyvers et al., 2004) to elicit insight. In this paper, we focus our experiments on patent data.

1.2 Primer on Latent Dirichlet allocation

Fortunately, the intuition behind LDA topic models is relatively straightforward. To understand how the algorithm infers the topics in relation to

documents we first define what constitutes a topic. A topic is a distribution over a fixed vocabulary, where each word has an assigned probability of occurrence. Subsequently, we can take the view that each document is likely a product of one or more topics, a cocktail of themes as it were with different proportions of each ingredient.

Take for example the following document sampled from the August 2015 EPO Worldwide Patent Statistical Database (PATSTAT). The patent abstract contained in Figure 1.1 relates to a mechanism for stopping a water wheel. We have taken the liberty of highlighting a selection of words from a few of this document's prominent topics. Words like "pressure", "liquid", and "flow" belong to the **fluids/water** topic and are colored blue. While words relating to the **mechanisms** by which this fluid is directed such as "chamber", "valve", and "guide" are colored red. Finally, words such as "transmission", "speed", and "operated" belong to the topic associated with **signals** and are colored green.

PURPOSE: To stop a **water wheel** stably by closing guide vanes to an opening at which **water hammer** phenomena scarcely occur and suppressing the shake of the **guide vanes** caused by the **transmission** of **water pressure** when a main **valve** is closed next.

CONSTITUTION: If a running **water wheel** receives a stop instruction at **time t1**, **guide vanes** G are closed gradually, and at the same **time**, **water wheel** load is decreased, and at **time t2** when this **water wheel** becomes no-load, a paralleling breaker is opened and also a governor S is cut out. The opening of the **guide vanes** at no-load continues to be closed after that, but when a safety pin is broken at **time t3**, this breakage is detected and the governor is **operated** again, the sound **guide vanes** are opened to a specified opening near the no-load opening and fixed, the over **speed** of the **water wheel** is prevented and the opening and closing moment of the **guide vanes** is balanced.'

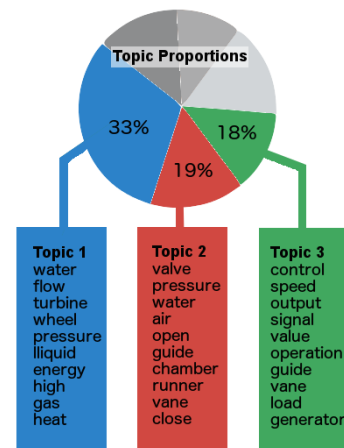


FIGURE 1.1: Topic proportions in a sample patent abstract.

The topics in the previous example were formed not over a single document but over a collection. The grey sections of the pie chart above represent the topics that this patent does not contain strong elements of. This is a key characteristic of LDA topic models, each document has a unique topic 'fingerprint' as a result of a generative process. That process for generating a document word by word is as follows. First we decide, sampling from the distribution of topics, which topic the word will belong to. Then we sample from that topic's distribution to decide what the word itself will be. This process is then simply repeated for each word, and while it works it has the following assumptions:

- Documents can manifest multiple topics (however typically not many)
- Each document is assumed to be the product of a generative process.
- Generative process starts with a topic, i.e. a distribution over a fixed vocabulary.
- Assumes a fixed number of topics

Latent Dirichlet allocation falls into a family of machine learning algorithms called **hidden variable models**. In this family of models, customarily the user "posits a hidden structure in the observed data, and then learns that structure using posterior probabilistic inference" (Blei and Lafferty, 2009). For LDA specifically, the documents are the observed data, the topics and document topic proportions are hidden.

More formally, we may define this process mathematically as a joint distribution over our hidden variables and our observed variables. Namely the distribution over vocabulary β , the topic proportions for document d θ_d , the topic assignment for a word in a document $T_{d,n}$ and of course the observed words themselves $w_{d,n}$.

$$P(\beta_{1:K}, \theta_{1:D}, T_{1:D}, W_{1:D}) = \prod_{i=1}^K p(\beta_i) \prod_{d=1}^D p(\theta_d) \left\{ \prod_{n=1}^N p(T_{d,n} | \theta_d) p(w_{d,n} | \beta_{1:K}, T_{d,n}) \right\} \quad (1.1)$$

In Eq 1.1 we see a few dependencies worth noting. Firstly that the topic we assign to a word $T_{d,n}$ depends on the distribution of topics of its document θ_d . Additionally, that the identity of the word itself is dependent on not only the topic we assigned to generate it $T_{d,n}$, but also the vocabulary distributions of each topic $\beta_{1:K}$. Equivalently we can express the dependencies between these variables as a graphical model, illustrated in figure 1.2.

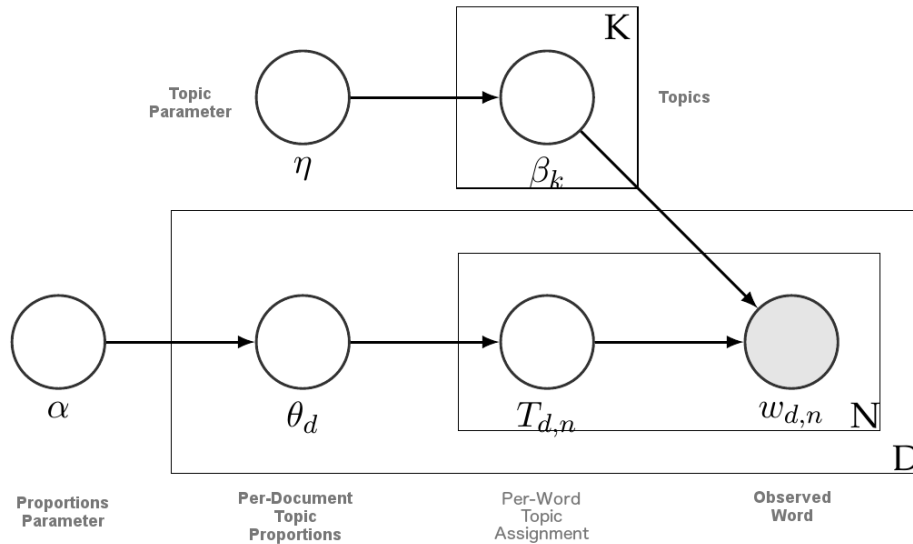


FIGURE 1.2: Graphical model for LDA

So how do we actually obtain our estimates of the hidden parameters? We need to calculate the conditional distribution of our hidden parameters (the topic structure), and the observed words i.e. the posterior distribution described in Eq. 1.2. However the denominator makes this calculation

computationally infeasible due to the number of combinations our hidden parameters could take.

$$p(\beta_{1:K}, \theta_{1:D}, Z_{1:D} | W_{1:D}) = \frac{p(\beta_{1:K}, \theta_{1:D}, Z_{1:D} | W_{1:D})}{p(W_{1:D})} \quad (1.2)$$

To move past this, most solutions use either sampling or variational based methods to perform approximate inference and obtain estimates of the hidden parameters. Variational methods allow us to translate the original problem to one of optimization and take advantage of the many optimization techniques available. This in turn allows us to make extensions that are often faster, scale better or allow for different forms of input such as streaming documents.

1.3 Adding a temporal component

One such extension, and the extension we explore in this study, is to relax the implicit assumption of LDA that the order of the documents doesn't matter. By incorporating the order of the documents to the model, a topic is no longer simply a distribution over words but now becomes a *sequence* of distributions over words. This is the jump that allows us not only to identify a theme, as with static LDA, but also track how it progresses in time, giving us the Dynamic Topic Model (DTM).

The DTM offers several advantages over traditional LDA including improved predictive performance (Blei and Lafferty, 2006). Primarily though, it facilitates a greater understanding of how each topic developed, and how the ideas therein formed and matured. With it, we can inspect trends of word usage to uncover a richer and more detailed hidden structure. For instance figure 1.3 contains a sample theme from a sub-collection of hydroelectric patents and the progression of word prevalences within it over time.

1.4 Why patents?

Patent data is specifically interesting in this context because of the role patents play in company formation, job growth, economic development, and novel invention. Their history tells a story of technological progression. In an attempt to maintain a competitive edge, many companies large and small spend a considerable amount of energy researching this history to identify technical trends relevant to their industry.

Dynamic topic models have the potential to aid this research by enabling us to track the progression of innovation through language use in patent abstracts. In this paper we look at a number of elements, including

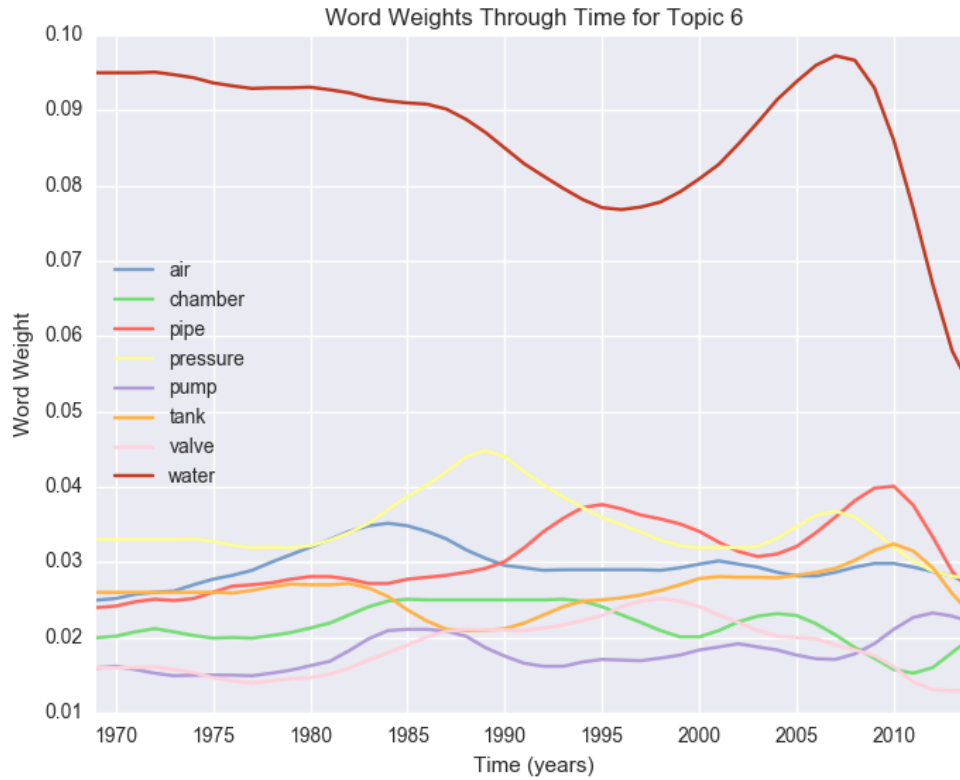


FIGURE 1.3: Distribution over words in a sample hydro-electric topic over time

the evolution of technological themes and their proportions, the origination and progression of language, as well as document influences. Furthermore, the patent corpus and associated International Patent Classification (IPC) labels provide a platform for the comparison of various topic modeling algorithms.

1.5 Experiments

At the time of writing this, surprisingly little has been published exploring the effectiveness of both the DTM and the DIM. Much research has evaluated model quality solely based on the likelihood of held-out predictions, however this does not always translate to semantically meaningful topics (Chang et al., 2009). Additionally, the predictive performance of these models is adversely affected by longer time horizons due to an "increase in the rate of specialization in scientific language" (Blei and Lafferty, 2006). Acknowledging the room to explore alternative methods of model evaluation, we implemented the experiments listed below.

Historical Topic Trend Validation

The simplest, but also the most hands-on, method of evaluating the quality of topics produced by the DTM and DIM is simply to validate the inferred topic trends against industry history. For instance, if in the topic of water purification systems we observe a rise in the usage of words "2D materials" and "lattice membranes" around 2005, we might substantiate this by pointing out graphene's isolation the previous year in 2004.

1.5.1 Topic Coherence

In light of research suggesting that likelihoods and perplexity don't always correlate with human judgement on the interpretability of topics (Blei and Lafferty, 2006) we borrow several methods of topic coherence suggested by (Rosner et al., 2014). Namely, we evaluated model topic coherences using C_v , C_{npmi} , C_{uci} , and U_{mass} . Using C_v , the metric most correlated with human judgement, DTM achieved the highest with score with XXX compared to static LDA with YYY. For complete results see Table ZZZ in Section 5.

1.5.2 Classification

Unsupervisedly learned word vectors have demonstrated useful as features in a myriad of NLP tasks. We wished to evaluate the proficiency of the word vectors generated by the DTM and DIM at creating an effective feature space for text categorization. To do this we made use of the IPC labels of patents as broad class labels for text content. The resulting topic vectors should then help identify which class a document belongs to. Naturally we tested the efficacy of each model's vector space at correctly classifying the IPC label of documents when fed to a range of classification algorithms. Peak classification performance of the DTM based classifiers was F1 XXX, while LDA was F1 YYY. Text classification results are given in section ZZZ.

1.5.3 Clustering

Another method by which we hoped to evaluate the quality of the resulting document vectors was by their ability to define separations in the data relative to the ground truth CPC labels. In order to assess which models yielded vector spaces of the corpus that most effectively clustered the documents we used the following metrics: the adjusted rand index, normalized mutual score info, homogeneity, completeness and the V-measure. Indeed we found that the DTM's vector space tended to outperform that of LDA at clustering with a peak NMI score of XXX compared to YYY. For more detailed results, refer to table ZZZ in section 5.

Chapter 2

Background Information and Theory

2.1 Literature Review

2.2 DTM model

2.2.1 Issues

addresses several latent structures in the document collection such as topic evolution and prevalence, however does not address the birth and death of topics, like models such as Ahmed and Xing, [2012](#).

2.3 DIM model

Chapter 3

Experimental Set Up

3.1 Data Prep and Considerations

Chapter 4

Experimental Results

4.1 DTM Results and Insights

4.1.1 Topics Through Time

validating topic histories in technology

4.2 DIM Results and Insights

4.2.1 Influence Metric

validating influential patents

correlation with forward citations

correlation with page-rank

4.3 Performance Evaluation

4.3.1 Classification

4.3.2 Clustering

Chapter 5

Usefulness in Other Models

5.1 Economic Model

influence can be used as a proxy for forward citations when citations are not available. used as a gamma in model for likelihood of innovation

Chapter 6

Conclusions and Future Work

6.1 Conclusions

6.2 Future Work

Appendix A

Appendix Title Here

Write your Appendix content here.

Bibliography

- Ahmed, Amr and Eric P. Xing (2012). "Timeline: A Dynamic Hierarchical Dirichlet Process Model for Recovering Birth/Death and Evolution of Topics in Text Stream". In: *CoRR* abs/1203.3463. URL: <http://arxiv.org/abs/1203.3463>.
- Blei, D. and J. Lafferty (2009). "Topic Models". In: *Text Mining: Theory and Applications*, Taylor and Francis.
- Blei, David M. and John D. Lafferty (2006). "Dynamic Topic Models". In: *Proceedings of the 23rd International Conference on Machine Learning*. ICML '06. Pittsburgh, Pennsylvania, USA: ACM, pp. 113–120. ISBN: 1-59593-383-2. DOI: [10.1145/1143844.1143859](https://doi.org/10.1145/1143844.1143859). URL: <http://doi.acm.org/10.1145/1143844.1143859>.
- Chang, Jonathan et al. (2009). "Reading Tea Leaves: How Humans Interpret Topic Models". In: *Neural Information Processing Systems*. Vancouver, BC. URL: [docs/nips2009-rt1.pdf](https://docs.nips2009-rt1.pdf).
- Newman, David et al. (2006). "Analyzing Entities and Topics in News Articles Using Statistical Topic Models". In: *Proceedings of the 4th IEEE International Conference on Intelligence and Security Informatics*. ISI'06. San Diego, CA: Springer-Verlag, pp. 93–104. ISBN: 3-540-34478-0, 978-3-540-34478-0. DOI: [10.1007/11760146_9](https://doi.org/10.1007/11760146_9). URL: http://dx.doi.org/10.1007/11760146_9.
- Rosner, Frank et al. (2014). "Evaluating topic coherence measures". In: *CoRR* abs/1403.6397. URL: <http://arxiv.org/abs/1403.6397>.
- Steyvers, Mark et al. (2004). "Probabilistic Author-topic Models for Information Discovery". In: *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD '04. Seattle, WA, USA: ACM, pp. 306–315. ISBN: 1-58113-888-1. DOI: [10.1145/1014052.1014087](https://doi.org/10.1145/1014052.1014087). URL: <http://doi.acm.org/10.1145/1014052.1014087>.