

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 23, 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	1
1.1 The need for topic models	1
1.2 Primer on Latent Dirichlet allocation	2
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
1.6 Roadmap of Paper	7
2 Background Information and Theory	8
2.1 Literature Review	8
2.1.1 Applications of Topic Modeling	8
2.1.2 Ensuring Model Quality	8
Perplexity Testing	8
Coherence Testing	9
2.1.3 Document Classification	10
Accuracy	10
Precision	10
Recall	11
F1 Score	11
2.1.4 Document Clustering	11
Adjusted Rand Score	11
Normalized Mutual Info	12
2.2 DTM Model Overview	13
2.2.1 Chaining models together	13
2.2.2 Variational approximate inference	14
3 Experimental Set Up	16
3.1 Considerations	16
3.2 Data Pre-Processing	17
3.3 Tuning models	17
3.4 Classification Set Up	19
3.5 Clustering Set Up	19

4	Experimental Results	20
4.1	DTM Results and Insights	20
4.1.1	Topics Through Time	20
	validating topic histories in technology	20
4.2	DIM Results and Insights	20
4.2.1	Influence Metric	20
	validating influential patents	20
	correlation with forward citations	20
	correlation with page-rank	20
4.3	Performance Evaluation	20
4.3.1	Classification	20
4.3.2	Clustering	20
5	Usefulness in Other Models	21
5.1	Economic Model	21
6	Conclusions and Future Work	22
6.1	Conclusions	22
6.2	Future Work	22
A	Appendix Title Here	23
	Bibliography	24

List of Figures

1.1	Patent114	2
1.2	Graphical model for LDA	3
1.3	wwtfTopic6	5
2.1	Graphical model for DTM showing a series of chained static LDA models. The triangles represent the Kalman filter esti- mates of the hyper-parameters	14
3.1	PSQLSchema	17
3.2	DTMCV	18

List of Tables

3.1	The effects of treatments X and Y on the four groups studied.	18
-----	---	----

Chapter 1

Introduction

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, Yahoo, and Bidu that help 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 derived 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 extracting topics from patent data with the hopes identifying meaningful trends in renewable energy technologies.

1.2 Primer on Latent Dirichlet allocation

Fortunately, the intuition behind LDA topic models is relatively straight forward. 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.'

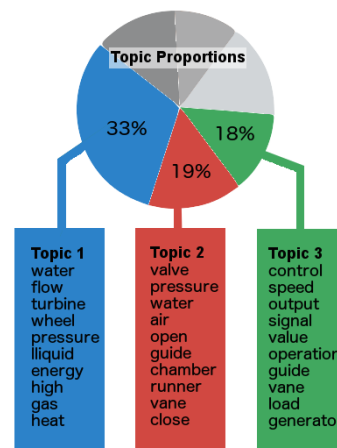


FIGURE 1.1: Topic proportions in a sample patent abstract.

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 our first word will belong to. Then we sample from that chosen 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 worth noting:

- 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, the user customarily "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. Specifically, we define the distribution over vocabulary as β , the topic proportions for document d θ_d , the topic assignment for a word in a document $z_{d,n}$ and of course the observed words themselves $w_{d,n}$.

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

In Eq 1.1 we see a few dependencies worth noting. Firstly that the topic we assign to a word $z_{d,n}$ 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 $z_{d,n}$, 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.

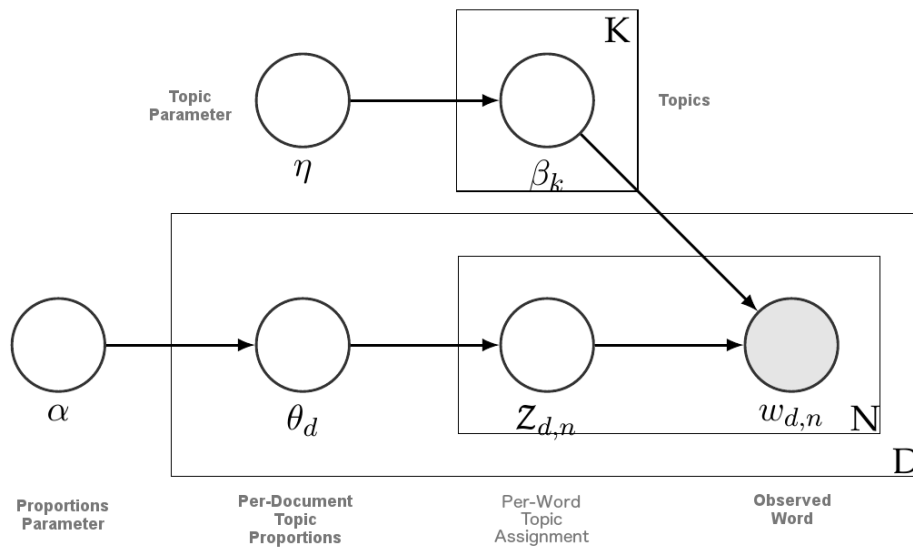


FIGURE 1.2: Graphical model for LDA

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.

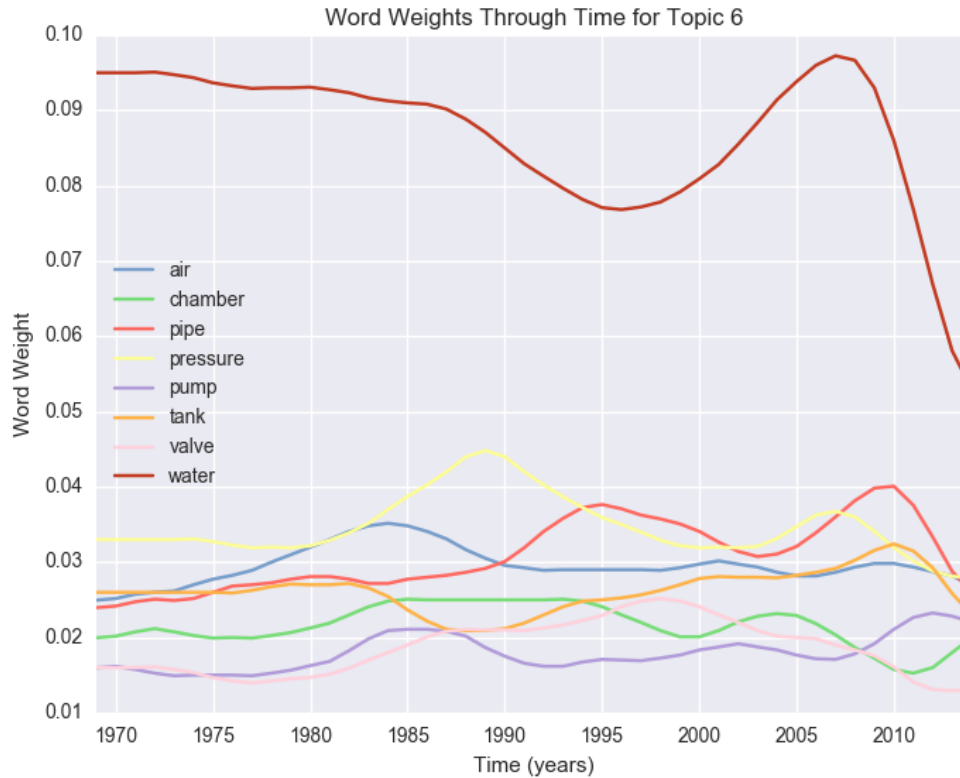


FIGURE 1.3: Distribution over words in a sample hydro-electric topic over time

Dynamic topic models have the potential to aid this research by enabling us to track the evolution of innovation through language use in patent abstracts. In this paper we look at a number of elements, including the evolution of technological themes and their proportions, the origination and development 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 likelihood 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 known 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

We wished to evaluate the proficiency of the word vectors unsupervisedly generated by the DTM and DIM at forming an effective feature space for document 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 we used to evaluate the quality of the resulting document vectors was by their ability to cluster the documents. In order to determine which models yielded vector spaces of the corpus that most effectively defined separations in the data relative to the ground truth CPC labels we used the following metrics: the adjusted rand index, normalized mutual score info, homogeneity, completeness and the V-measure which we cover in section 2.1. 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.

1.6 Roadmap of Paper

The overall structure of this paper is as follows. In **chapter 2** we review the literature surrounding the applications and various evaluation methods for topic models, and also give a detailed account of both the DTM and DIM. Then in **Chapter 3** we cover the experimental set up and considerations in data preparation. The results of our experiments are subsequently presented in **Chapter 4**, and finally we conclude in **Chapter 5** with a discussion of results and suggestions for future study.

Chapter 2

Background Information and Theory

2.1 Literature Review

In this section we begin by reviewing a few of the tasks common in topic modeling. Then we describe a handful of the ways the quality of topic models, LDA in particular, are commonly tested. Finally we discuss the specific tasks of document classification and clustering in the context of topic modeling, as well as the corresponding methods for evaluating model performance at these tasks.

2.1.1 Applications of Topic Modeling

The most popular application of topic models is simply summarizing large text collections by mining the topics. This is a task LDA is particularly suited for (Griffiths and Steyvers, 2004; Mei, Shen, and Zhai, 2007). The original LDA paper however (Blei, Ng, and Jordan, 2003) gave promising results on document classification as well. Since then LDA has been used with success not only for document classification, but also for clustering and information retrieval (Wei and Croft, 2006; Nagwani, 2015). This is due to the strength of the topic vectors LDA models provide, which tend to correlate strongly with human judgement.

2.1.2 Ensuring Model Quality

Perplexity Testing

In order to ensure the strength of these topic vectors researchers employ a handful methods to evaluate the topic models. While the most intuitive method is simply to have humans judge the coherence of each topic, this becomes prohibitively time consuming and expensive for large data sets. One commonly used method of automating this process is by evaluating the topic model on a held out set of testing documents and obtaining the log-likelihood perplexity of the unseen documents (Blei, Ng, and Jordan,

2003; Wallach et al., 2009). A higher likelihood on unseen documents, and a lower perplexity score indicates a better model. However this method of evaluating topic model performance has several issues. Firstly, it has been shown that predictive likelihood, or equivalently perplexity, is not always correlated with human judgement, and in some cases is even slightly anti correlated (Chang et al., 2009). Secondly this method of evaluation only acts as a general measure of the entire model. What about the quality of the individual topics?

Coherence Testing

Fortunately several methods of evaluating the coherence of individual topics from topic models exist. For a topic t we define the **Umass** coherence as a sum of the pairwise scores of that topic's top words $W_t = \{w_1, \dots, w_n\}$.

$$\begin{aligned} \text{Umass Coherence } c(t, W_t) &= \sum_{w_i, w_j \in W_t} \text{score}(w_i, w_j) \\ &= \sum_{w_i, w_j \in W_t} \log \frac{d(w_i, w_j) + \epsilon}{d(w_i)} \end{aligned} \quad (2.1)$$

Where $d(w_i)$ is the number of documents containing the word w_i and $d(w_i, w_j)$ is the number of documents containing both word w_i and w_j . The ϵ in the numerator is simply to smooth the counts and is typically set to a minimal value such as 1 or .01. Intuitively then, a topic is good if its words cooccur often (Mimno et al., 2011).

The **UCI** measure introduced by (Newman, Bonilla, and Buntine, 2011), operates in the same manner as Umass but with the pointwise mutual information as a scoring function instead, given in eq 2.2.

$$\begin{aligned} \text{UCI Coherence } = c(t, W_t) &= \sum_{w_i, w_j \in W_t} \text{score}(w_i, w_j) \\ &= \sum_{w_1, w_2 \in W_t} \log \frac{p(w_i, w_j)}{p(w_i)p(w_j)} \end{aligned} \quad (2.2)$$

Where $p(w_i)$ is the probability of seeing word w_i in a random document and $p(w_i, w_j)$ is the probability of seeing both word w_i and word w_j together in a random document. It should be noted that obtaining these probabilities requires empirically estimating them from an external dataset.

Two more noteworthy measures of topic coherence, in addition to those outlined above, were developed by Roder, Both and Hinneberg in their study titled "Exploring the Space of Topic Coherence Measures" (Röder, Both, and Hinneburg, 2015). These measures were the **C_v** and **C_{npmi}**

measures which demonstrated a substantial correlation with human judgement. For brevity we do not replicate their derivations here, but the interested reader will find a detailed description of each in (Röder, Both, and Hinneburg, 2015)

2.1.3 Document Classification

Though the topics produced by topic models are useful in their own right for the qualitative analysis of documents, they are also useful quantitatively when trying to classify documents. For instance a large news organization may want to automatically sort its thousands of articles into the categories "politics", "natural disasters" and "sports". To do this they might use a topic model to get a vector of topic proportions for each document to use as features for a classification algorithm. This process is referred to as document vectorization.

While baseline methods for document vectorization exist, such as the Term Frequency Inverse Document Frequency (tf-idf), LDA has been shown to outperform them in certain scenarios. For instance when less training data is available LDA boasted a shorter training time and higher classification accuracy (Li and Zhang, 2010). Additionally when tested against other baseline methods for document vectorization such as the unigram model or probabilistic latent semantic analysis (PLSA), LDA again proved consistently more accurate at document classification tasks (Lu, Mei, and Zhai, 2011).

Accuracy

When it comes time to evaluate a model's classification performance there are several approaches, the most intuitive of which is accuracy. Accuracy is defined as

$$acc = \frac{Tp + Tn}{Tp + Tn + Fp + Fn} \quad (2.3)$$

Where Tp are our true positives, Tn are our true negatives, Fp are our false positives, and Fn are our false negatives. It should be noted that normal values of accuracy for classification tasks depend highly on the data at hand. Noisy data or a large number of classes can both artificially drive accuracy scores down.

Precision

But what if we want to know, out of the total number of guesses for a particular class, what fraction were correct? For this, researchers typically use

the Precision, defined as

$$p = \frac{Tp}{Tp + Fp} \quad (2.4)$$

Recall

Conversely, if we wish to know out of the total number of cases we could have guessed correctly, what fraction we *did* guess correctly then Recall is typically used. Recall is defined as

$$r = \frac{Tp}{Tp + Fn} \quad (2.5)$$

F1 Score

The F1 score is a way of combining the above two metrics Precision and Recall into one wholistic measure. Conceptually it is the harmonic mean of the Precision and Recall where we assign even weights to each. The F1 score is defined as

$$\begin{aligned} F_1 &= \frac{1}{\frac{1}{2} \left(\frac{1}{p} + \frac{1}{r} \right)} \\ &= \frac{2pr}{p + r} \end{aligned} \quad (2.6)$$

2.1.4 Document Clustering

Another well established task for topic models is document clustering. LDA has been used to successfully cluster a range of documents such as news articles and legal judgements (Lu, Mei, and Zhai, 2011; Xie and Xing, 2013; Kumar and Raghuveer, 2013). As opposed to classification where we want to assign an explicit label to each document, with clustering we wish to evaluate how well the resulting document topic vectors separate the documents into a meaningful structure.

Adjusted Rand Score

One way of accomplishing this is by using the **Adjusted Rand Score** which measures the similarity of two sets of class labels; namely the true labels C and those predicted by a clustering algorithm K . We may calculate the raw (unadjusted) Rand index following equation 2.7 (Hubert and Arabie, 1985).

$$RI = \frac{a + b}{C_2^{n_{\text{samples}}}} \quad (2.7)$$

Where a is the number of pairs of elements in C belonging to the same class, and in K belonging to the same class. Conversely b is the number of pairs of elements in C belonging to different classes, and in k belonging to different classes. Finally, $C_2^{n_{\text{samples}}}$ is the total number of possible pairs in the dataset. In order to ensure that random labelings receive a score of zero we define the adjusted Rand index as

$$ARI = \frac{RI - E[RI]}{\max(RI) - E[RI]} \quad (2.8)$$

Normalized Mutual Info

The Normalized Mutual Info is another method of evaluating clustering performance that has been successfully applied in the context of topic modelling (Xu, Liu, and Gong, 2003; Cai et al., 2008). It again assumes we have two sets of labels, this time we call them U and V , over N objects. We define the entropy of a label set U in equation 2.9, where $P(i) = |U_i|/N$ is the probability that a random object from U falls into class U_i .

$$H(U) = \sum_{i=1}^{|U|} P(i) \log(P(i)) \quad (2.9)$$

The mutual information (MI) between U and V can be expressed as

$$MI(U, V) = \sum_{i=1}^{|U|} \sum_{j=1}^{|V|} P(i, j) \log\left(\frac{P(i, j)}{P(i)P(j)}\right) \quad (2.10)$$

With these two components we can write the normalized mutual information as proposed in (Vinh, Epps, and Bailey, 2009).

$$NMI(U, V) = \frac{MI(U, V)}{\sqrt{H(U)H(V)}} \quad (2.11)$$

2.2 DTM Model Overview

This section briefly outlines the dynamic topic model (**DTM**), following closely the original derivation found in (Blei and Lafferty, 2006). As this is intended as more of a summary, we recommend the reader examine the original paper for a complete exposition of the mechanics of the DTM.

In our primer on LDA (in section 1.2) we outlined the conceptual basis for static LDA topic models. Namely, that topics consist of a distribution over a fixed vocabulary and are determined by a set of hyper-parameters β . Additionally each document is represented as a combination of topics, with proportions controlled by their corresponding set of hyper-parameters α . Roughly speaking, the goal of the Dynamic Topic Model (**DTM**) is to account for the drift in topics over time by chaining together a series of static LDA models. This is accomplished by tying the hyper-parameters $\alpha_{t-1}, \beta_{t-1}$ at time step $t - 1$, to the hyper-parameters α_t, β_t at time step t . The result is a model that allows us to track how our topics evolve at each time step.

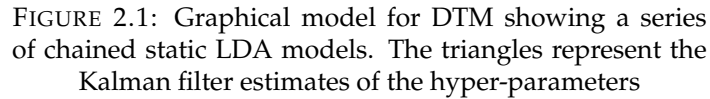
2.2.1 Chaining models together

The question then, is how do we tie our hyper-parameters together? Well, regularly with static LDA we would simply use a Dirichlet distribution to model our uncertainties in word distributions (hence the name). Unfortunately, the Dirichlet distribution does not lend itself to sequential modeling, which eliminates this option. Instead, we make a state-space model that evolves with Gaussian noise to chain together the natural parameters of each topic $\beta_{z,t}$ such that each topic "evolves" from the last.

$$\beta_{z,t} | \beta_{z,t-1} \sim \mathcal{N}(\beta_{t-1,z}, \sigma^2 I) \quad (2.12)$$

Similarly, with static LDA we would also pull our document specific topic proportions θ from the Dirichlet distribution. For the same reason as above this is no longer an option. So to express our uncertainty over our topic proportions, we use a logistic normal with mean α . Then we chain our topic proportions together using the same trick as we did above with word distributions, (by using Gaussian noise). This yields the graphical model in figure 2.1.

$$\alpha_t | \alpha_{t-1} \sim \mathcal{N}(\alpha_{t-1}, \delta^2 I) \quad (2.13)$$



Because we have used the Gaussian distribution to model the progression of our parameters, inference becomes intractable due to the non-conjugacy of Gaussian and multinomial models. To get around this we take the same variational approach to approximate inference as before in section 1.2 with static LDA. Taking a variational approach has the advantage of allowing us to handle larger document sets compared to Gibbs sampling which becomes computationally difficult at large corpus sizes.

We begin by creating a collection of variational parameters we will optimize over our latent variables. Our latent variables are the topics $\beta_{t,k}$, topic

proportions $\theta_{t,d}$, and topic indicators $Z_{t,d,n}$. While we have variational parameters for each topic (consisting of a sequence of multinomial parameters), and for each document (the latent topic proportions). The resulting posterior, again following the notation of (Blei and Lafferty, 2006) is given by equation 2.14.

$$\prod_{k=1}^K q(\beta_{k,1}, \dots, \beta_{k,T} | \hat{\beta}_{k,1}, \dots, \hat{\beta}_{k,T}) \times \prod_{t=1}^T \left(\prod_{d=1}^{D_t} q(\theta_{t,d} | \gamma_{t,d}) \prod_{n=1}^{N_{t,d}} q(z_{t,d,n} | \phi_{t,d,n}) \right) \quad (2.14)$$

This is where we tune our approximate posterior, and specifically the variational observations $\{\hat{\beta}_{k,1}, \dots, \hat{\beta}_{k,T}\}$ according to the KL divergence between the estimated and the true posterior. Note that here, each topic proportions vector $\gamma_{t,d}$ receives a corresponding free Dirichlet parameter, while each topic indicator $z_{t,d,n}$ receives a corresponding free multinomial parameter $\phi_{t,d,n}$. To optimize the document topic proportion vectors we subsequently employ gradient ascent, however this is not necessary for the document level parameter updates as they simply have a closed form.

Finally, we may track our variational parameters $\hat{\beta}$ and $\hat{\alpha}$ between time slices using either the Kalman filter or wavelet regression. For brevity we will not replicate the mechanics of these methods here and encourage the interested reader to refer to the derivations provided in detail in (Blei and Lafferty, 2006).

Chapter 3

Experimental Set Up

In the interest of reproducibility we include in this section some of the considerations specific to our data set. Additionally we cover the data pre-processing steps taken prior to our experiments and describe our process for model tuning. Finally, we outline the procedure taken for the classification and clustering experiments.

3.1 Considerations

The data we used for our experiments comes from the August 2015 EPO Worldwide Patent Statistical Database (**PATSTAT**). According to the EPO, this data set contains over 90 million patent documents from leading industrialized and developing countries, with some documents as early as the nineteenth century. Thankfully, the expanse of this data is met with an equal level of documentation¹.

Due to time and computation limits, we ran experiments only on a subset of this data rather than the whole corpus. In order to break the PATSTAT into manageably sized subclasses we made use of Cooperative Patent Classification (**CPC**) labels. The CPC labels are a hierarchical system used globally to annotate patents according to the area of technology to which they relate. Predominantly we carried out our research on patents under the umbrella of **Y02E 10/20**, the subclass of patents relating to **hydro energy**.

An additional consideration in filtering our data was the International Patent Documentation Centre (**INPADOC**) **patent family**. A patent family is a collection of patents filed in various countries to protect the same invention. We ensure that only one member of each patent family participates in the study so as not to 'double count' the same invention. Furthermore, only patents filed in English are considered. Figure 3.1 contains the schema for the PostgreSQL database we constructed.

¹For more information the official catalogue can be found at <http://goo.gl/LRQWnu>

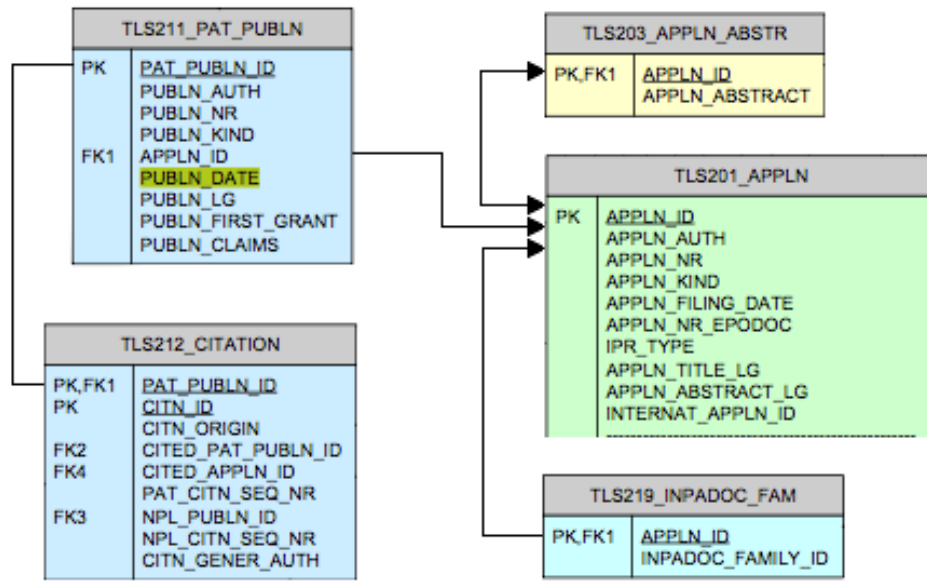


FIGURE 3.1: Schema of PostgreSQL database containing relevant tables

3.2 Data Pre-Processing

The phrase "garbage in garbage out" is commonly used in machine learning to emphasize the importance of the data pre-processing step. We begin our pre-processing by retrieving abstracts from the aforementioned PSQL database and applying case folding such that words like "Turbine" at the beginning of a sentence will match the word "turbine" elsewhere. We then remove unnecessary symbols and punctuation via regex and apply stopping using the English **NLTK** stopwords list. This effectively removes common conjunctions and operational words from our corpus that don't actually contribute to the topics, such as "but, if, the, and" etc.

We then further distill the abstracts by applying lemmatization. This step is meant to reduce inflectional forms of a word such as "operates, operating, operational" to a base form "operate". After lemmatization, we remove any words occurring less than 25 times total, matching the implementation found in (Blei and Lafferty, 2006). Finally, we save the resulting corpus, totaling over 6,400 documents, and serialize the dictionary of vocabulary for subsequent access.

3.3 Tuning models

Our selection of hyperparameters can have a large influence on the topics that models infer. For this reason, prior to running any of our experiments, we must take care in tuning the hyperparameters of our topic models, namely the number of topics K and the maximum number of iterations

TABLE 3.1: The effects of treatments X and Y on the four groups studied.

Groups	Treatment X	Treatment Y
1	0.2	0.8
2	0.17	0.7
3	0.24	0.75
4	0.68	0.3

M_{iter} .

To do this we defined a range of values for these parameters which we then sampled parameter sets from uniformly. At each parameter set we evaluated the umass, uci, npmi, and cv topic coherence of the resulting model. Table XXX contains the best parameter sets found for each model according to each coherence measure on a subset of documents. Figure 3.2 displays the cv coherences of the DTM as a function of K and M_{iter} . From this plot we can see that increasing the number of topics K beyond a certain point fails to improve topic coherence with similar results for the number of iterations.

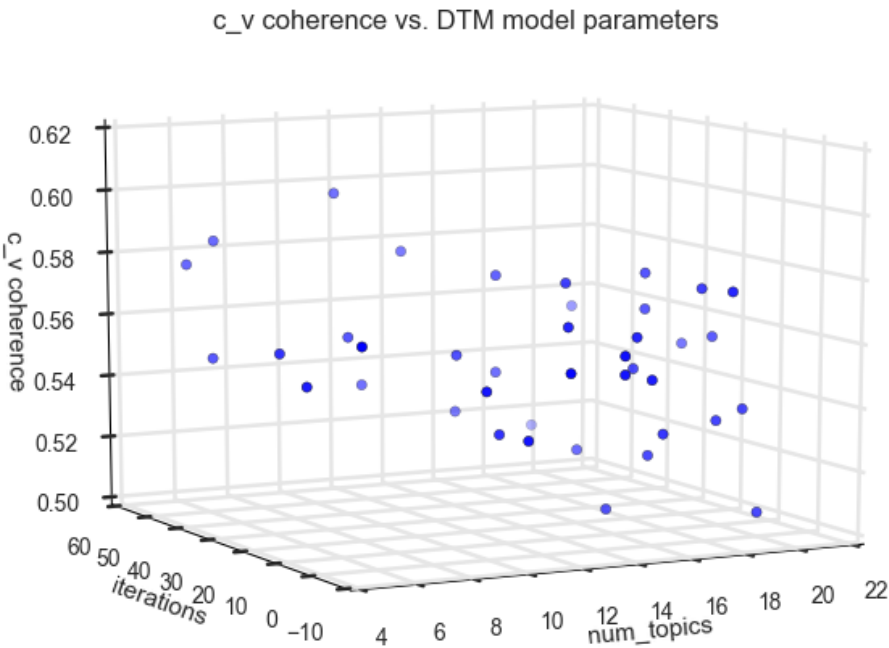


FIGURE 3.2: cv coherence values for DTM as a function of K and M_{iter}

you can see from the plot that (some comment or observation here.)

Ultimately went with cv recommended parameter sets as it correlated the highest with human judgement in (reference here)

setting the alpha parameter in LDA as recommended by (e.g., as in Steyvers and Griffiths 2007) It acts as a form of regularization. A recommended setting is $\alpha = 50/K$ (Griffiths and Steyvers 2004; Steyvers and Griffiths 2007)

3.4 Classification Set Up

3.5 Clustering Set Up

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

The clustering result is evaluated by comparing the Normalized mutual information (Xu et al. 2003; Cai et al. 2008) [NEED TO ACTUALLY CITE]

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

- 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>.
- Blei, David M., Andrew Y. Ng, and Michael I. Jordan (2003). "Latent Dirichlet Allocation". In: *J. Mach. Learn. Res.* 3, pp. 993–1022. ISSN: 1532-4435. URL: <http://dl.acm.org/citation.cfm?id=944919.944937>.
- Cai, Deng et al. (2008). "Modeling Hidden Topics on Document Manifold". In: *Proceedings of the 17th ACM Conference on Information and Knowledge Management*. CIKM '08. Napa Valley, California, USA: ACM, pp. 911–920. ISBN: 978-1-59593-991-3. DOI: [10.1145/1458082.1458202](https://doi.org/10.1145/1458082.1458202). URL: <http://doi.acm.org/10.1145/1458082.1458202>.
- 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).
- Griffiths, T. L. and M. Steyvers (2004). "Finding scientific topics". In: *Proceedings of the National Academy of Sciences* 101.Suppl. 1, pp. 5228–5235.
- Hubert, Lawrence and Phipps Arabie (1985). "Comparing partitions". In: *Journal of Classification* 2.1, pp. 193–218. ISSN: 1432-1343. DOI: [10.1007/BF01908075](https://doi.org/10.1007/BF01908075). URL: <http://dx.doi.org/10.1007/BF01908075>.
- Kumar, R. V. and K. Raghuvver (2013). "Legal Documents Clustering and Summarization using Hierarchical Latent Dirichlet Allocation". In: *IAES International Journal of Artificial Intelligence* 1.2, pp. 27–35. URL: <http://dx.doi.org/10.11591/ij-ai.v2i1.1186>.
- Li, Lei and Yimeng Zhang (2010). "An empirical study of text classification using Latent Dirichlet Allocation". In:
- Lu, Yue, Qiaozhu Mei, and Chengxiang Zhai (2011). "Investigating Task Performance of Probabilistic Topic Models: An Empirical Study of PLSA and LDA". In: *Inf. Retr.* 14.2, pp. 178–203. ISSN: 1386-4564. DOI: [10.1007/s10791-010-9141-9](https://doi.org/10.1007/s10791-010-9141-9). URL: <http://dx.doi.org/10.1007/s10791-010-9141-9>.
- Mei, Qiaozhu, Xuehua Shen, and Chengxiang Zhai (2007). "Automatic Labeling of Multinomial Topic Models". In: *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD '07. San Jose, California, USA: ACM, pp. 490–499. ISBN: 978-1-59593-609-7. DOI: [10.1145/1281192.1281246](https://doi.org/10.1145/1281192.1281246). URL: <http://doi.acm.org/10.1145/1281192.1281246>.

- Mimno, David et al. (2011). "Optimizing Semantic Coherence in Topic Models". In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. EMNLP '11. Edinburgh, United Kingdom: Association for Computational Linguistics, pp. 262–272. ISBN: 978-1-937284-11-4. URL: <http://dl.acm.org/citation.cfm?id=2145432.2145462>.
- Nagwani, N. K. (2015). "Summarizing large text collection using topic modeling and clustering based on MapReduce framework". In: *Journal of Big Data* 2.1, pp. 1–18. ISSN: 2196-1115. DOI: [10.1186/s40537-015-0020-5](https://doi.org/10.1186/s40537-015-0020-5). URL: <http://dx.doi.org/10.1186/s40537-015-0020-5>.
- Newman, David, Edwin V. Bonilla, and Wray L. Buntine (2011). "Improving Topic Coherence with Regularized Topic Models." In: *NIPS*. Ed. by John Shawe-Taylor et al., pp. 496–504. URL: <http://dblp.uni-trier.de/db/conf/nips/nips2011.html#NewmanBB11>.
- 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.
- Röder, Michael, Andreas Both, and Alexander Hinneburg (2015). "Exploring the Space of Topic Coherence Measures". In: *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*. WSDM '15. Shanghai, China: ACM, pp. 399–408. ISBN: 978-1-4503-3317-7. DOI: [10.1145/2684822.2685324](https://doi.org/10.1145/2684822.2685324). URL: <http://doi.acm.org/10.1145/2684822.2685324>.
- 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>.
- Vinh, Nguyen Xuan, Julien Epps, and James Bailey (2009). "Information Theoretic Measures for Clusterings Comparison: Is a Correction for Chance Necessary?" In: *Proceedings of the 26th Annual International Conference on Machine Learning*. ICML '09. Montreal, Quebec, Canada: ACM, pp. 1073–1080. ISBN: 978-1-60558-516-1. DOI: [10.1145/1553374.1553511](https://doi.org/10.1145/1553374.1553511). URL: <http://doi.acm.org/10.1145/1553374.1553511>.
- Wallach, Hanna M. et al. (2009). "Evaluation Methods for Topic Models". In: *Proceedings of the 26th International Conference on Machine Learning (ICML)*. Ed. by Léon Bottou and Michael Littman. Montreal: Omnipress, pp. 1105–1112.
- Wei, Xing and W. Bruce Croft (2006). "LDA-based Document Models for Ad-hoc Retrieval". In: *Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. SIGIR '06. Seattle, Washington, USA: ACM, pp. 178–185. ISBN: 1-59593-369-7. DOI: [10.1145/1148170.1148204](https://doi.org/10.1145/1148170.1148204). URL: <http://doi.acm.org/10.1145/1148170.1148204>.

- Xie, Pengtao and Eric P. Xing (2013). "Integrating Document Clustering and Topic Modeling". In: *CoRR* abs/1309.6874. URL: <http://arxiv.org/abs/1309.6874>.
- Xu, Wei, Xin Liu, and Yihong Gong (2003). "Document Clustering Based on Non-negative Matrix Factorization". In: *Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. SIGIR '03. Toronto, Canada: ACM, pp. 267–273. ISBN: 1-58113-646-3. DOI: [10.1145/860435.860485](https://doi.org/10.1145/860435.860485). URL: <http://doi.acm.org/10.1145/860435.860485>.