
A Language-based Approach to Measuring Scholarly Impact

Sean M. Gerrish

David M. Blei

Department of Computer Science, Princeton University 35 Olden St., Princeton, NJ 08540

SGERRISH@CS.PRINCETON.EDU

BLEI@CS.PRINCETON.EDU

Abstract

Identifying the most influential documents in a corpus is an important problem in many fields, from information science and historiography to text summarization and news aggregation. Unfortunately, traditional bibliometrics such as citations are often not available. We propose using changes in the thematic content of documents over time to measure the importance of individual documents within the collection. We describe a dynamic topic model for both quantifying and qualifying the impact of these documents. We validate the model by analyzing three large corpora of scientific articles. Our measurement of a document's impact correlates significantly with its number of citations.

1 Introduction

Measuring the influence of a scientific article is an important and challenging problem. Influence measurements are used to assess the quality of academic instruments, such as journals, scientists, and universities, and can play a role in decisions surrounding publishing and funding. They are also important for academic research: finding and reading the influential articles of a field is central to good research practice.

The traditional method of assessing an article's influence is to count the citations to it. The impact factor of a journal, for example, is based on aggregate citation counts (Garfield, 2002). This is intuitive: if more people have cited an article, then more people have read it, and it is likely to have had more impact on its field. Citation counts are used with other types of documents as well: the Pagerank algorithm, which uses hyperlinks of web-pages, has been essential to Google's early success in Web search (Brin & Page, 1998).

Appearing in *Proceedings of the 27th International Conference on Machine Learning*, Haifa, Israel, 2010. Copyright 2010 by the author(s)/owner(s).

Though citation counts can be powerful, they can be hard to use in practice. Some collections, such as news stories, blog posts, or legal documents, contain articles that were influential on others but lack explicit citations between them. Other collections, like OCR scans of historical scientific literature, do contain citations, but they are difficult to read in reliable electronic form. Finally, citation counts only capture one kind of influence. All citations from an article are counted equally in an impact factor, when some articles of a bibliography might have influenced the authors more than others.

We take a different approach to identifying influential articles in a collection. Our idea is that an influential article will affect how future articles are written and that this effect can be detected by examining the way corpus statistics change over time. We encode this intuition in a time-series model of sequential document collections.

We base our model on dynamic topic models, allowing for multiple threads of influence within a corpus (Blei & Lafferty, 2006). Though our algorithm aims to capture something different from citation, we validate the inferred influence measurements by comparing them to citation counts. We analyzed one hundred years of the Proceedings of the National Academy, one hundred years of *Nature*, and a forty year corpus of articles on computational linguistics. With only the language of the articles as input, our algorithm produces a meaningful measure of each document's influence in the corpus.

2 The Document Influence Model

We develop a probabilistic model that captures how past articles exhibit varying influence on future articles. Our hypothesis is that an article's influence on the future is corroborated by how the language of its field changes subsequent to its publication. In the model, the influence of each article is encoded as a hidden variable and posterior inference reveals the influential articles of the collection.

Our model is based on the dynamic topic model (DTM) (Blei & Lafferty, 2006), a model of sequential corpora that allows language statistics to drift over time. Pre-

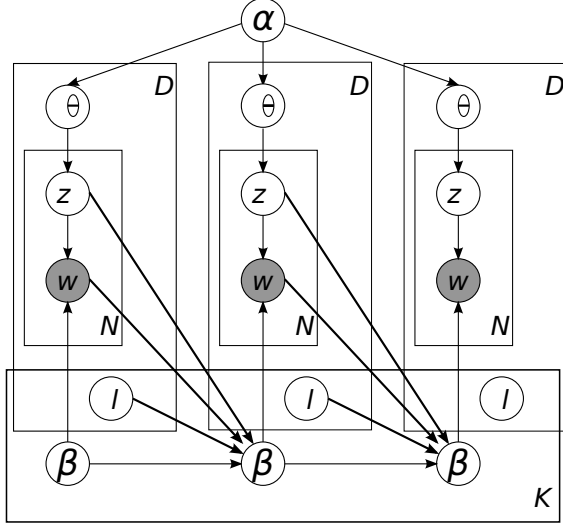


Figure 1. The Document Influence Model.

vious probabilistic models of text assumed that the underlying distributions over words were fixed. The DTM introduced a Markov chain of term distributions to capture probabilities that drift over the course of the collection.

Let V be the number of words in a vocabulary and consider the natural parameters β_t of a term distribution at time t , where the probability of a word w is given by the softmax transformation of the unconstrained vector,

$$p(w | \beta_t) \propto \exp(\beta_{t,w}). \quad (1)$$

The corresponding distribution over terms, i.e., the “topic,” is a point on the vocabulary simplex. In the logistic normal Markov chain, this distribution drifts with the stationary autoregressive process

$$\beta_{t+1} | \beta_t \sim \mathcal{N}(\beta_t, \sigma^2 I), \quad (2)$$

where σ^2 is the transition variance.

Now consider a corpus with D articles at each time t and let the rows $\mathbf{w}_{t,d}$ of $\mathbf{W}_{t,1:D}$ denote the articles as vectors of word counts. In the simplest DTM, a single distribution over words drifts according to Equation 2. For each time point, the words of its articles are drawn independently from Equation 1. One can then compute the posterior distribution of the sequence of topics $\beta_{1:T}$ conditioned on the observed documents. This summarizes the corpus as a smooth trajectory of word frequencies.

We now turn to our central idea: some articles influence the topic more than others. In our model, each article is assigned a normally distributed *influence score* ℓ_d , which is a scalar value that describes the influence that article d has on the topic. The higher the influence, the more the words of the article affect how the topic drifts.

This is encoded in the time series model. The more influential a document is, the more its words “nudge” the topic’s natural parameters at the next time step,

$$\beta_{t+1} | \beta_t, (w, \ell)_{t,1:D} \sim \mathcal{N}(\beta_t + \exp(-\beta_t) \sum_d w_{d,t} \ell_{t,d}, \sigma^2 I). \quad (3)$$

The words of an article with a high influence will have a higher expected probability in the next epoch; the words of an article with zero influence will not affect the next epoch.

We call this model the *document influence model* (DIM). Conditioned on a corpus, the posterior distribution of the topic and influence scores gives a trajectory of term frequencies and a retrospective estimate of the influence of each article. An article whose words can help explain the way the word frequencies change will have a high posterior influence score. We show in Section 4.2 that this estimate of influence is meaningful.

Multiple topics. Corpora typically contain multiple persistent themes. Accordingly, the full dynamic topic model contains multiple topics, each associated with a time series of distributions. Conditioned on the topics, articles at each time are modeled with latent Dirichlet allocation (LDA). Each article exhibits the topics with different random proportions θ_d ; each word of each article is drawn by choosing a topic assignment from those proportions $z_{d,n}$, and choosing a word from the corresponding topic (Blei et al., 2003).

Modeling multiple topics is important to the influence model because an article might have different impact in the different fields that it discusses. For example, an article about computational genomics may be very important to biology but less important to computer science. We want to discern its influence on each of these topics separately.

As with the DTM, we posit K topic trajectories, and each document of each time point is modeled with LDA. For each document, we now associate an influence score $\ell_{d,k}$ for each topic k . Each of the K topics drifts according to an adapted version of Equation 2, where we restrict attention to the influence score for that topic and to the words of each document that were assigned to it,

$$\beta_{k,t+1} | \beta_{k,t}, (w, \ell, z)_{t,1:D} \sim \mathcal{N}(\beta_{k,t} + \exp(-\beta_{k,t}) \sum_d \ell_{d,k} \sum_n w_{d,n} z_{d,n,k}, \sigma^2 I). \quad (4)$$

Here, $z_{d,n,k}$ is the indicator that the n th word in document d is assigned to topic k and we have dropped the index t on z and w . The graphical model is illustrated Figure 1.

Although we presented our model in this section with influence spanning one year, we also adapted it to accommodate an “influence envelope”, where an article’s influence spans W years. This provides a more realistic model of influence

(Porter et al., 1988), but it complicates the inference algorithm and may not be necessary, as we note in section 4.2.

To use this model, we analyze a corpus through posterior inference. This reveals a set of K changing topics and influence scores for each article and each topic. The posterior provides a thematic window into the corpus and can help identify which articles most contributed to the development of its themes.

Related work. There is a large literature on citation analysis and bibliometrics. See Osareh (1996) for a review. Much work in this area uses the link structure of citation networks to extract higher level structure. Borner et al. (2003) for example, have used author and citation networks to understand the evolution of ideas in the history of science. There has also been work that proposes features and models for predicting future citation counts. Successful features often include the publishing journal’s impact factor, previous citations to last author, key terms, and number of authors (Tang & Zhang, 2009; Lokker et al., 2008).

A number of algorithms include the text of the articles in their analysis. This work often models the information in citations by predicting them or modeling them with topics (Nallapati & Cohen, 2008; Chang & Blei, 2009; Dietz et al., 2007; Cohn & Hofmann, 2001) or other semantic tools (McNee et al., 2002; Ibáñez et al., 2009). Other work in this area uses the text of documents along with citations to summarize documents (Qazvinian & Radev, 2008) or to propose new bibliometrics: Mann et al. (2006) use topic models and citations to map topics over time and define several new bibliometric measurements such as topic Impact Factor, topical diffusion, and topic longevity.

Our model has a different flavor from this research. We are interested in identifying the influential articles in a collection, but we do not assume that there are any notions of reference within them: there is no training data that contains citations. While we validate our model by looking at the relationship between our measure of influence and citation counts, our model is applicable to collections for which this kind of information does not exist.

Two important pieces of recent research have similar goals. Leskovec et al. (2009) describe a framework for tracking the spread of memes, or ideas, in document collections, and investigate the direction in which ideas tend to percolate. Shaparenko & Joachims (2007) describe a measure of influence by modeling documents as unigram mixtures of earlier documents and use a likelihood ratio test to predict citations between documents. In contrast to this work, the DIM uses dynamic topics to explicitly model the change in *topic* language. Further, we do not attempt to model links between documents, as in Shaparenko & Joachims (2007).

3 Inference and parameter estimation

Our computational challenge is to compute the posterior distribution of the latent variables—the sequences of topics and the per-document influence values—conditioned on an observed corpus. As for simpler topic models, this posterior is intractable to compute exactly. We employ variational methods to approximate it. Variational methods posit a simpler distribution over the latent variables with free parameters (called variational parameters). These parameters are fit to minimize the KL divergence between the variational distribution and the true posterior, which is equivalent to maximizing a lower bound on the marginal probability of the observations. See Jordan et al. (1999) for a review of this class of approximate inference methods.

We begin by specifying a variational distribution for the DIM posterior. First, the word assignments z_n and topic proportions θ_d are governed by multinomial parameters ϕ_d and Dirichlet parameters γ_d , as in LDA (Blei et al., 2003).

The variational distribution for topic trajectories $\{\beta_{k,1}, \dots, \beta_{k,T}\}$ is described by a linear Gaussian chain. It is governed by parameters $\{\tilde{\beta}_{k,1}, \dots, \tilde{\beta}_{k,T}\}$, which are interpreted as the “variational observations” of the chain. These induce a sequence of means \tilde{m}_t and variances \tilde{V}_t . Blei & Lafferty (2006) call this a “variational Kalman filter.”

Finally, the variational distribution of the document influence value $\ell_{d,k}$ is a Gaussian with mean $\tilde{\ell}_{d,k}$ and fixed variance σ_{ℓ}^2 .

The variational distribution is

$$q(\beta, \ell, z, \theta | \tilde{\beta}, \tilde{\ell}, \phi, \gamma) = \prod_{k=1}^K q(\beta_{k,1:T} | \tilde{\beta}_{k,1:T}) \prod_{t=1}^T \prod_{d=1}^{D_t} q(\theta_{t,d} | \gamma_{t,d}) q(\ell_d | \tilde{\ell}_d) \prod_{n=1}^{N_{t,d}} q(z_{t,d,n} | \phi_{t,d,n}).$$

Using this variational family, our goal is to maximize the following lower bound on the model evidence of the observed words \mathbf{W} :

$$\begin{aligned} \ln p(\mathbf{W}) &\geq \sum_T \mathbb{E}_q [\ln p(\beta_t | \beta_{t-1})] \\ &+ \sum_T \sum_{D_t} \mathbb{E}_q [\ln p(\ell_d)] + \mathbb{E}_q [\ln p(\theta_d | \alpha)] \\ &+ \sum_T \sum_{D_t} \sum_{N_d} \mathbb{E}_q [\ln p(z_n | \theta_d)] + \mathbb{E}_q [\ln p(w_n | z_n, \beta_t)] \\ &+ H(q). \end{aligned}$$

This bound is optimized by coordinate ascent, with the variational parameters optimized sequentially in blocks. These updates are repeated until the relative increase in the lower bound is below a threshold.

Topic trajectories. The variational update for $\tilde{\beta}$ is similar to that in Blei & Lafferty (2006). For each topic, we

update the variational Kalman “observations” by applying gradient ascent:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \beta_{sw}} = & -\frac{1}{\sigma^2} \sum_{t=1}^T (\tilde{m}_{tw} - \tilde{m}_{t-1,w} - G_{t-1,w}) \\ & \times \left(\frac{\partial \tilde{m}_{tw}}{\partial \beta_{sw}} - \frac{\partial \tilde{m}_{t-1,w}}{\partial \beta_{sw}} + G_{t-1,w} \frac{\partial \tilde{m}_{t-1,w}}{\partial \beta_{sw}} \right) \\ & + \sum_T \left(N_{w,t} - N_t \zeta_t^{-1} \exp(\hat{m}_{\beta_{tw}} + \frac{\tilde{V}_{tw}}{2}) \right) \frac{\partial \tilde{m}_{tw}}{\partial \beta_{sw}} \\ & + \frac{1}{\sigma^2} \sum_{t=1}^T \frac{\partial \tilde{m}_{t-1,w}}{\partial \beta_{sw}} (H_{t-1,w} - G_{t-1,w}^2) \\ & + \frac{1}{\sigma^2} \sum_{t=0}^{T-1} \frac{\partial \tilde{m}_{tw}}{\partial \beta_{sw}} G_{tw} \tilde{V}_{tw}, \end{aligned}$$

where

$$\begin{aligned} G_{sn} &= \mathbb{E}_q [\exp(-\beta_{s,k,n}) (\mathbf{W}_{s,k,n} \circ z_{s,k,n}) \ell_{s,k}] \\ H_{sn} &= \mathbb{E}_q [\exp(-2\beta_{s,k,n}) ((\mathbf{W}_{s,k,n} \circ z_{s,k,n}) \ell_{s,k})^2]. \end{aligned}$$

We expand H_t in the supplementary materials. Note also the variational parameter ζ_t and the term $\frac{\partial \tilde{m}_{tw}}{\partial \beta_{sw}}$, both described in Blei & Lafferty (2006). The former can be updated once per iteration with $\zeta_t \leftarrow \sum_w \exp(\tilde{m}_{tw} + \tilde{V}_{t,n}/2)$. The latter can be derived from the variational Kalman filter updates (see the supplementary materials).

Influence values. In the DIM, changes in a topic’s mean parameters are governed by a normal distribution. As a consequence of this choice, updates for the influence parameters $\tilde{\ell}_{t,k}$ solve a linear regression. In this regression, documents’ words at time t explain the expected topic drift $\Delta_{\beta,t,k} = (\beta_{t+1,k} - \beta_{t,k})$, where the contributions of each document’s words are given by the design matrix $X = \text{Diag}(\exp(-\beta_{t,k})) (\mathbf{W}_{t,k} \circ \phi_{t,k})$. ($\text{Diag}(\vec{x})$ refers to the matrix having the elements of \vec{x} on its diagonal.)

The parameter updates for document influence $\tilde{\ell}_{t,k}$ are defined, for each time t and each topic k , by the variational normal equation

$$\tilde{\ell}_{t,k} \leftarrow \left(\frac{\sigma^2}{d} \mathbf{I} + \mathbb{E}_q [X^T X] \right)^{-1} \mathbb{E}_q [X^T \Delta_{\beta,t,k}].$$

The expectation $\mathbb{E}_q [X^T X]$ is a matrix with dimension $D_t \times D_t$. Its elements are

$$\begin{aligned} \mathbb{E}_q [X^T X]_{d,d'} = & \sum_n \exp(-2\tilde{m}_{t,k,n} + 2\tilde{V}_{t,k,n}) \\ & \times (w_{t,d,n} w_{t,d',n} \phi_{t,k,d,n} \phi_{t,k,d',n}) \end{aligned}$$

when $d \neq d'$ and

$$\begin{aligned} \mathbb{E}_q [X^T X]_{d,d} = & \sum_n \exp(-2\tilde{m}_{t,k,n} + 2\tilde{V}_{t,k,n}) \\ & \times (w_{t,d,n}^2 \phi_{t,k,d,n}) \end{aligned}$$

otherwise. The expectation $\mathbb{E}_q [X^T \Delta_{\beta,t,k}]$ is a D_t -

dimensional matrix with elements

$$\begin{aligned} \mathbb{E}_q [X^T \Delta_{\beta,t,k}]_d = & \sum_n w_{t,d,n} \phi_{t,k,d,n} \\ & \times (\tilde{m}_{t+1,k,n} - \tilde{m}_{t,k,n} + \tilde{V}_{t,k,n}/2) \\ & \times \exp(-\tilde{m}_{t,k,n} + \tilde{V}_{t,k,n}/2). \end{aligned}$$

Topic proportions and topic assignments. Updates for the variational Dirichlet on the topic proportions $\theta_{d,k}$ have a closed-form solution, exactly as in LDA (Blei et al., 2003); we omit details here.

The variational parameter for each word w_n ’s hidden topic z_n is the multinomial ϕ_n . We solve for $\phi_{n,k}$ by the closed-form updates

$$\begin{aligned} \log(\phi_{n,k}) \leftarrow & \Psi(\gamma_k) + \tilde{m}_{t,k,n} \\ & + \frac{1}{\sigma^2} w_t \tilde{\ell}_{d_n,k} \exp(-\tilde{m}_{t,k} + \tilde{V}_{t,k}/2) (\tilde{m}_{t+1,k} - \tilde{m}_{t,k} + \tilde{V}_{t,k}) \\ & - \frac{1}{\sigma^2} w_{t,n} \left[\tilde{\ell}_{d_n,k} \exp(-2\tilde{m}_{t,k} + 2\tilde{V}_{t,k}) \right. \\ & \quad \left. \times (\mathbf{W}_{t,n,\setminus d_n} \circ \phi_{t,n,k,\setminus d_n}) \tilde{\ell}_{t,k,\setminus d_n} \right] \\ & - \frac{1}{\sigma^2} w_{t,n}^2 \exp(-2\tilde{m}_{t,k} + 2\tilde{V}_{t,k}) (\tilde{\ell}_{d_n,k}^2 + \sigma_l^2) \end{aligned}$$

where Ψ is the digamma function. Solving the constrained optimization problem, this update is followed by normalization, $\phi_{w,k} \leftarrow \frac{\phi_{w,k}}{\sum_K \phi_{n,k}}$.

4 Empirical study

We studied three sequential corpora of scientific articles. For each corpus, we estimated and examined the posterior distributions of its articles’ influence.

In this section, we demonstrate that the estimate of an article’s influence is robustly correlated to the number of citations it received. While the DIM model is designed for corpora without citations—and, indeed, only the documents’ text and dates are used in fitting the model—citations remain an established measure of influence. This study provides validation of the DIM as an exploratory tool of influential articles.

4.1 Data

The three corpora we analyzed were the *ACL Anthology*, *The Proceedings of the National Academy of Science*, and the journal *Nature*. For each corpus, we removed short documents, terms that occurred in too few documents, and terms that occurred in too many documents. We also removed terms whose statistics did not vary over the course of the collection, as such terms would not be useful for assessing change in language (a sample of such non-varying terms from *Nature* is “ordinarily”, “shake”, “centimetre”, “traffic”, and “themselves”).

ACL Anthology. The *Association for Computational Linguistics Anthology* is a digital collection of publications about computational linguistics and natural language processing (Bird et al., 2008). We analyzed a 50% sample from this anthology, spanning 1964 to 2002. Our sample contains 7,561 articles and 11,763 unique terms after preprocessing. For this corpus we used article citation counts from the *ACL Anthology Network* (Radev et al., 2009).

PNAS. The *Proceedings of the National Academy of Sciences* is a leading, highly-cited, multidisciplinary scientific journal covering biological, physical, and social sciences. We sampled one seventh of the collection, spanning 1914 (when it was founded) to 2004. Our sample contains 12,145 articles and 14,504 distinct terms after preprocessing. We found citations using Google Scholar for 78% of this collection.

Nature. The journal *Nature* is the world’s most highly cited interdisciplinary science journal (Thompson Reuters) with content on a range of scientific fields. We analyzed a 10% sample from this corpus, spanning 1869 (when it was founded) to 2008. Our sample contains 34,418 articles and 6,125 distinct terms after preprocessing. We found citations using Google Scholar for 31% of these documents.

Inference for 10 topics on each corpus above took about 11 hours to converge on a desktop Intel 2.4GHz Core 2 Quad cpu. Our convergence criterion was met when the evidence lower bound increased by no more than 0.01%. For the experiments described below, we set topics’ Markov chain variance $\sigma^2 = 0.005$ and $\sigma_d = \sigma_l = 0.0001$.

4.2 Relating posterior influence and citation

We studied the DIM with varying numbers of topics. We measured the relationship between the posterior influence values of each article $\tilde{\ell}_d$ and its citation count c_d .

We first aggregate the influence values across topics. Recall that each document has an influence value for each topic. For each word, we compute its expected posterior influence score, with the expectation taken with respect to its (random) topic assignment. We then sum these values over all words in the document,

$$f(\tilde{\ell}_d) = \sum_{n=1}^{N_d} E[z_{d,n} \cdot \tilde{\ell}_d]. \quad (5)$$

This weights each word by the influence associated with its assigned topic. (Using the maximum value of influence across topics yielded similar results.)

Figure 2 displays the Spearman rank correlation between the aggregated posterior influence score of Equation 5 and citation counts. The DIM posterior—which is estimated only from the texts of the articles—has a positive correlation to the number of citations. All of these numbers were

found significant up to $p < 1e - 4$, using permutation tests on the influence scores.

Correlation goes up when we model multiple topics within a corpus. Moving from 2 to 5 topics in the *ACL* corpus increases correlation from 0.25 to 0.37. *Nature* is likewise better with more topics, with a correlation of 0.28 at 20 topics; while *PNAS* performs best near 5 topics, with a correlation of 0.20.

Figure 2 also shows the fraction of citations explained by DIM scores: *Nature* documents with the highest 20% of posterior influence, for example, received 56% of citations. The flat regions in *ACL* and *PNAS* are due to aggregate influence scores very close to zero.

Heuristic model. The DIM is a complicated model. To justify its complexity, we describe a simple baseline (the *heuristic*) which captures our intuition with a single topic, is easy to implement, and runs quickly. For this heuristic, we define a word’s weight at time t as:

$$w_t := \frac{\text{Frequency of } w \text{ in } [t, t+f]}{\text{Frequency of } w \text{ in } [t-p, t]},$$

for fixed distances f into the future and p into the past. A document’s score is the weighted average of its words’ weights. This heuristic captures the intuition that influential documents use language adopted by other documents.

The heuristic performed best with large values of its parameters ($f = p = 200$). With these settings, it achieves a correlation of 0.20 for the *ACL*, 0.20 for *PNAS*, and 0.26 for *Nature*. For *Nature*, the model is more correlated with citations than the heuristic for 20, 50, and 75 topics. Correlation is matched for *PNAS*, the model slightly beating the heuristic at 5 topics. *ACL* outperforms the heuristic for all numbers of topics.

Shuffled corpus Though we have eliminated date as a confounder by controlling for it in correlations, there may be other confounders such as document length or topic distribution. We therefore measured the DIM’s relationship to citations when dates were randomly shuffled, keeping all documents which share a date together. If non-date confounders exist, then we might see correlation in the shuffled data, marking observed correlation as dubious.

We shuffled dates in the corpora and refit the DIM. We found a *maximum* date-controlled correlation of 0.018 for 29 shuffles of *ACL*; 0.001 for 5 shuffles of *Nature*; and 0.012 for 28 shuffles of *PNAS*. While this shuffled experiment and controlling for date do not entirely preclude confounding, they eliminate many potential confounders.

4.3 A closer look

Experiments showing correlation with citations demonstrate consistency with existing bibliometrics. However,

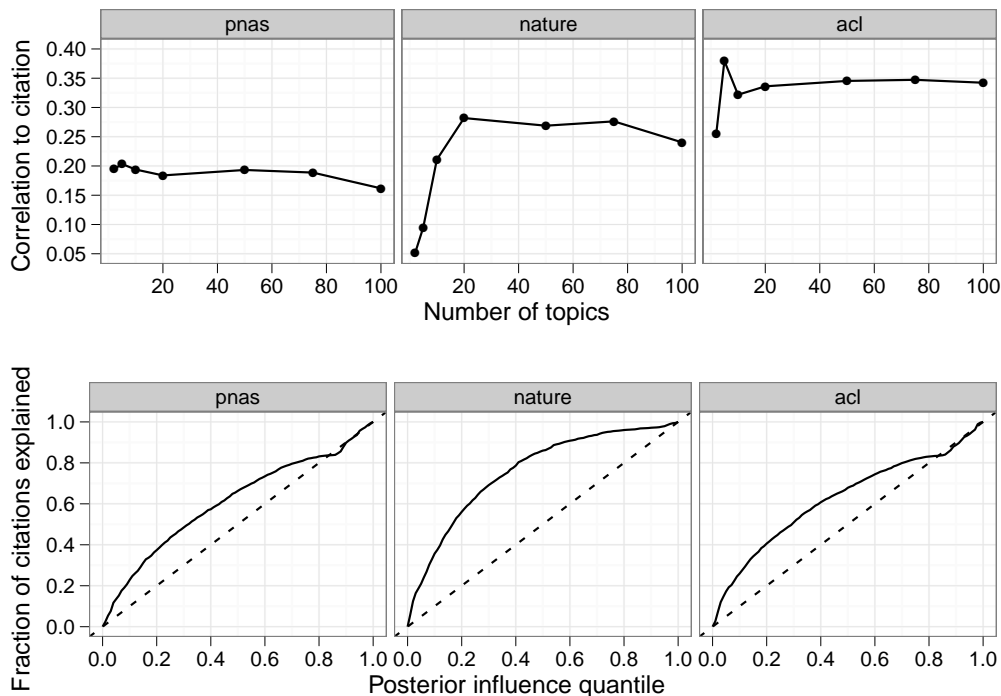


Figure 2. Spearman rank correlation between citation counts and posterior influence score, controlling for date (top) and fraction of citations explained by posterior influence (bottom).

the DIM also finds qualitatively different articles than citation counts. In this section we describe several documents to give the reader an intuition behind the kind of analysis that the DIM provides.

IBM Model 3 The second-most cited article in the *ACL Anthology Network* is *The Mathematics of Statistical Machine Translation: Parameter Estimation* (Brown et al., 1993). It has 450 intra-*ACL* citations and 2,130 total citations listed on Google Scholar. This seminal work describes parameter estimation for five word-based statistical models of machine translation; it provided widely accepted statistical models for word alignment and introduced the well-known “IBM models” for machine translation. The posterior influence score for Brown et al. (1993) ranked 6 out of 7,561 articles in a 10-topic model.

This article was most influential in a topic about translation, which had a trend toward “alignment for machine translation.” The largest-moving words are shown in Figure 3 (left). Upward trends for “alignment”, “brown”, and “equation” are evident (although it is not clear whether “brown” refers to the author or the corpus).

The Penn Treebank The most-cited article in our subset of the *ACL Anthology Network* is *Building a large annotated corpus of English: the Penn Treebank* (Marcus et al., 1993), with 1,622 *ACL* citations and 2,810 citations on Google Scholar. This article describes the large-scale part-of-speech and syntax tagging of a 4.5-million word corpus.

It falls in a topic about part-of-speech tagging and syntax trees; “treebank” had become one of the top words in the topic by 2004.

The DIM assigned a relatively low influence score to this article, ranking it 2,569 out of 7,561 articles. While Marcus et al. (1993) introduces a powerful *resource*, most of the article uses conventional language and ideas to detail the annotation of the Penn Treebank. As such, the paper does not discuss paradigm-changing ideas and the model scores it low. We emphasize that this does not undermine the tremendous influence that the Penn Treebank has had on the field of natural language processing. The DIM is not designed to discover this kind of influence.

Success in 1972 In 1967, The College Science Improvement Program was established to assist predominantly undergraduate institutions. Two years later *Nature* published a short column, which has the highest of our posterior influence in a 20-topic model, out of 34,418 *Nature* articles. No citation information was available about this article in Google Scholar. The column, *How to be Overtaken by Success*, discusses a debate about the “Miller bill”, which considers funding for postgraduate education (*Nature*). *Overtaken by Success* provides few research resources to researchers, which may explain lack of citation information. Instead, it presciently discusses a paradigm shift in a topic about science, industry, research, and education: “The record of the hearings [on the bill] is not merely an indication of the way the wind is blowing but an important

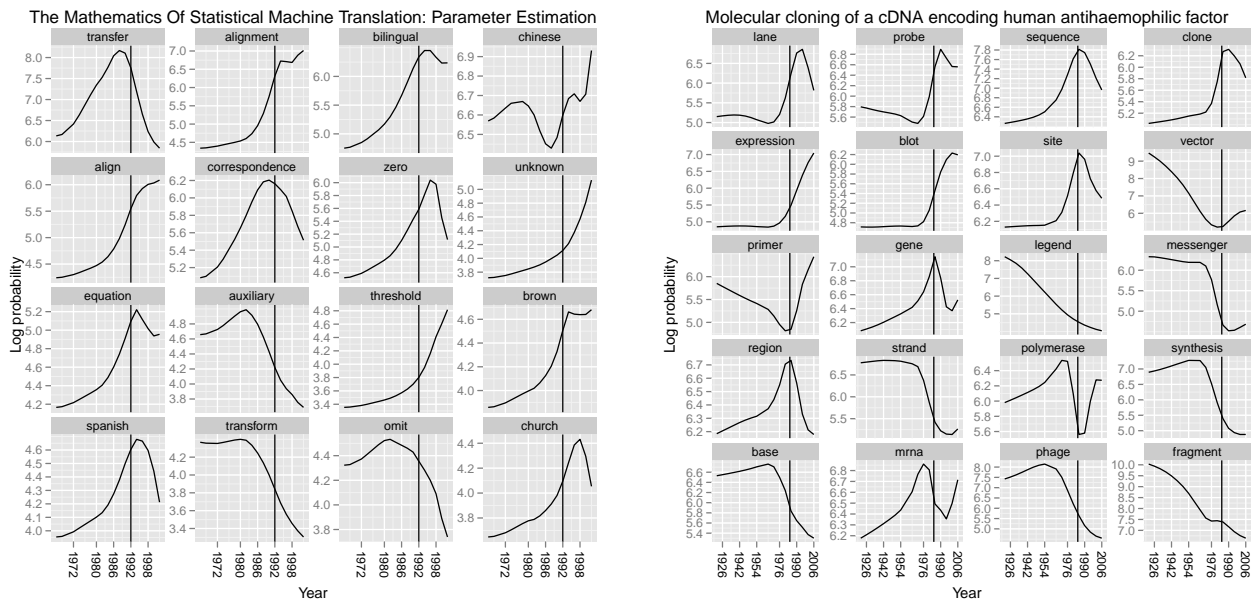


Figure 3. Most active words appearing in (Brown et al., 1993) (left) which have changed the most in a topic about translation. On right are words appearing in (Toole et al., 1984) in a topic about DNA and genetics. Terms are sorted by increase over 10 years.

guide to some of the strains which are now accumulating within the system of higher education...”

In 1972, three years after this article’s publication, The NSF Authorization Act of 1973 made the NSF explicitly responsible for science education programs *at all levels* (NSF, 2010). Where this may have been missed by those using citation counts to study the history of science education, the DIM has provided a metric with which to gauge interest in the article.

Genetics in Nature The sixth most influential document by the DIM in a 20-topic model of *Nature* is *Molecular cloning of a cDNA encoding human antihemophilic factor*, an article describing successful cloning of a human mRNA sequence important in blood clotting (Toole et al., 1984). With 584 citations, this article is among the top 0.2% of these 34,418 documents. The most active words appearing in this article are shown in Figure 3 (right). The plot shows some of the document’s key words – “expression”, “primer”, “blot” – become prominent words in the topic.

5 Conclusions and future work

We presented a time-series topic model for finding influential documents in long running sequential corpora. Based only on the changing statistics of the language in a corpus, we computed a measure of influence that is significantly related to observed citation counts. It would be useful to better understand how this metric is qualitatively different from citations and other bibliometrics: expert judgment or usage information obtained from digital libraries might be

some avenues. We leave this for future work.

The DIM could be made more realistic and more powerful in many ways. In one variant, individual documents might have their own “windows” of influence. Other improvements may change the way ideas themselves are represented, e.g. as atomic units, or *memes* (Leskovec et al., 2009). Further variants might differently model the flow of ideas, by modeling topics as birth and death processes, using latent force models (Alvarez et al., 2009), or by tracking influence *between* documents, building on the ideas of Shaparenko & Joachims (2007) or Dietz et al. (2007).

Acknowledgements

We would like to thank the reviewers for their insightful comments and JSTOR, Inc. for access to *PNAS*. David M. Blei is supported by a Google Research Grant, ONR 175-6343, and NSF CAREER 0745520.

References

A timeline of NSF history, 2010. URL <http://www.nsf.gov/about/history/>. [Online; accessed 31-January-2010].

Alvarez, Mauricio, Luengo, David, and Lawrence, Neil D. Latent force models. 2009.

Bird, Stephen, Dale, Robert, Dorr, Bonnie J., Gibson, Bryan, Joseph, Mark T., Kan, Min-Yen, Lee, Dongwon, Powley, Brett, Radev, Dragomir R., and Tan, Yee Fan. The ACL Anthology reference corpus: A reference dataset for bibliographic research. 2008.

- Blei, David and Lafferty, John. Dynamic topic models. *Proc. of the 23rd ICML*, 2006.
- Blei, David M., Ng, Andrew Y., and Jordan, Michael I. Latent Dirichlet allocation. *JMLR*, 2003.
- Borner, Katy, Chen, Chaomei, and Boyack, Kevin. Visualizing knowledge domains. *Annual Review of Information Science and Technology*, 37, 2003.
- Brin, Sergey and Page, Lawrence. The anatomy of a large-scale hypertextual web search engine. In *Computer Networks and ISDN Systems*, pp. 107–117, 1998.
- Brown, Peter F., Pietra, Vincent J., Della, Pietra, Stephen A. Della, and Mercer, Robert L. The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*, 19:263–311, 1993.
- Chang, Jonathan and Blei, David M. Relational topic models for document networks. *Proceedings of the 12th International Conference on Artificial Intelligence and Statistics (AISTATS) 2009*, 5, 2009.
- Cohn, David and Hofmann, Thomas. The missing link - a probabilistic model of document content and hypertext connectivity, 2001.
- Dietz, Laura, Bickel, Steffen, and Scheffer, Tobias. Unsupervised prediction of citation influences. *Proc. of the 24th ICML*, 2007.
- Garfield, Eugene. Algorithmic citation-linked historiography - mapping the literature of science. *ASIST 2002: Information, Connections, and Community*, 2002.
- Ibáñez, Alfonso, Larrañaga, Pedro, and Bielza, Concha. Predicting citation count of bioinformatics papers within four years of publication. *Bioinformatics*, 25:3303–3309, October 2009.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., and Saul, L. K. An introduction to variational methods for graphical models. *Learning in Graphical Models*, 1999.
- Leskovec, Jure, Backstrom, Lars, and Kleinberg, Jon. Meme-tracking and the dynamics of the news cycle. In *Proc. 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, New York, NY, USA, 2009. ACM.
- Lokker, Cynthia, McKibbin, K Ann, McKinlay, R. James, Wilczynski, Nancy L., and Haynes, R. Brian. Prediction of citation counts for clinical articles at two years using data available within three weeks of publication: retrospective cohort study. *The British Medical Journal*, March 2008.
- Mann, Gideon, Mimno, David, and McCallum, Andrew. Bibliometric impact measures leveraging topic analysis. June 2006.
- Marcus, Mitchell, Santorini, B., and Marcinkiewicz, M.A. Building a large annotated corpus of english: the Penn Treebank. *Computational Linguistics*, 19, 1993.
- McNee, Sean M., Albert, Istvan, Cosley, Dan, Gopalkrishnan, Prateep, Lam, Shyong K., Rashid, Al Mamunur, Konstan, Joseph A., and Riedl, John. On the recommending of citations for research papers. *Proc. 2002 ACM conference on Computer supported cooperative work*, pp. 116–125, 2002.
- Nallapati, Ramesh and Cohen, William. Link-plsa-lda: A new unsupervised model for topics and influence of blogs. *International Conference for Weblogs and Social Media*, 2008.
- Nature. How to be overtaken by success, April 1969.
- Osareh, Farideh. Bibliometrics, citation analysis and co-citation analysis: A review of literature i. *Libri*, 46:149–158, September 1996.
- Porter, A. L., Chubin, D. E., and Jin, Xiao-Yin. Citations and scientific progress: Comparing bibliometric measures with scientist judgments. *Scientometrics*, 13(3-4): 103–124, March 1988.
- Qazvinian, Vahed and Radev, Dragomir R. Scientific paper summarization using citation summary networks. pp. 689–696, 2008.
- Radev, Dragomir R., Joseph, Mark Thomas, Gibson, Bryan, and Muthukrishnan, Pradeep. A Bibliometric and Network Analysis of the field of Computational Linguistics. *Journal of the American Society for Information Science and Technology*, 2009.
- Shaparenko, B. and Joachims, T. Information genealogy: Uncovering the flow of ideas in non-hyperlinked document databases. 2007.
- Tang, Jie and Zhang, Jing. A discriminative approach to topic-based citation recommendation. *Advances in Knowledge Discovery and Data Mining*, 5476:572–579, 2009.
- Thompson Reuters. Journal Citation Reports - Science Edition, 2009.
- Toole, John J., Knopf, John L., Wozney, John M., Sultzman, Lisa A., Buecker, Janet L., Pittman, Debra D., Kaufman, Randal J., Brown, Eugene, Shoemaker, Charles, Orr, Elizabeth C., Amphlett, Godfrey W., Foster, W. Barry, Coe, Mary Lou, Knutson, Gaylord J., Fass, David N., and Hewick, Rodney M. *Nature*, 312:342–347, November 1984.