

Identification of Promising Research Directions using Machine Learning Aided Medical Literature Analysis

Victor Andrei and Ognjen Arandjelović
School of Computer Science
University of St Andrews
St Andrews KY16 9SX
Fife, Scotland
United Kingdom

Abstract—The rapidly expanding corpus of medical research literature presents major challenges in the understanding of previous work, the extraction of maximum information from collected data, and the identification of promising research directions. We present a case for the use of advanced machine learning techniques as an aide in this task and introduce a novel methodology that is shown to be capable of extracting meaningful information from large longitudinal corpora, and of tracking complex temporal changes within it.

I. INTRODUCTION

Recent years have witnessed a remarkable convergence of two broad trends. The first of these concerns information i.e. data – rapid technological advances coupled with an increased presence of computing in nearly every aspect of daily life, have for the first time made it possible to acquire and store massive amounts of highly diverse types of information. Concurrently and in no small part propelled by the environment just described, research in artificial intelligence – in machine learning [1], [2], [6], [3], data mining [11], and pattern recognition, in particular – has reached a sufficient level of methodological sophistication and maturity to process and analyse the collected data, with the aim of extracting novel and useful knowledge [5], [11]. Though it is undeniably wise to refrain from overly ambitious predictions regarding the type of knowledge which may be discovered in this manner, at the very least it is true that few domains of application of the aforesaid techniques hold as much promise and potential as that of medicine and health in general.

Large amounts of highly heterogeneous data types are pervasive in medicine. Usually the concept of so-called “big data” in medicine is associated with the analysis of Electronic Health Records [14], [4], [7], [22], [23], large scale sociodemographic surveys of death causes [19], social media mining for health related data [12] etc. Much less discussed and yet arguably no less important realm where the amount of information presents a challenge to the medical field is the medical literature corpus itself. Namely, considering the overarching and global importance of health (to say nothing of practical considerations such as the availability of funding), it is not surprising to observe that the amount of published medical research is immense and its growth is only continuing to accelerate. This presents a clear challenge to a researcher. Even restricted to a specified field of research, the amount of published data and findings makes it impossible for a human to survey the entirety of relevant publications exhaustively which inherently leads to the question as to what kind of important information or insight may go unnoticed or

insufficiently appreciated. The premise of the present work is that advanced machine learning techniques can be used to assist a human in the analysis of this data. Specifically, we introduce a novel methodology based on Bayesian non-parametric inference that achieves this, as well as free software which researchers can use in the analysis of their corpora of interest.

1) *Previous work*: A limitation of most models described in the existing literature lies in their assumption that the data corpus is static. Here the term ‘static’ is used to describe the lack of any associated temporal information associated with the documents in a corpus – the documents are said to be exchangeable [13]. However, research articles are added to the literature corpus in a temporal manner and their ordering has significance. Consequently the topic structure of the corpus changes over time [15], [9], [10]: new ideas emerge, old ideas are refined, novel discoveries result in multiple ideas being related to one another thereby forming more complex concepts or a single idea multifurcating into different ‘sub-ideas’ etc. The premise in the present work is that documents are not exchangeable at large temporal scales but can be considered to be at short time scales, thus allowing the corpus to be treated as *temporally locally static*.

II. PROPOSED APPROACH

In this section we introduce our main technical contributions. We begin by reviewing the relevant theory underlying Bayesian mixture models, and then explain how the proposed framework employs these for the extraction of information from temporally varying document corpora.

A. Bayesian mixture models

Mixture models are appropriate choices for the modelling of so-called heterogeneous data whereby heterogeneity is taken to mean that observable data is generated by more than one process (source). The key challenges lie in the lack of observability of the correspondence between specific data points and their sources, and the lack of *a priori* information on the number of sources [20].

Bayesian non-parametric methods place priors on the infinite-dimensional space of probability distributions and provide an elegant solution to the aforementioned modelling problems. Dirichlet Process (DP) in particular allows for the model to accommodate a potentially infinite number of mixture components [16]:

$$p(x|\pi_{1:\infty}, \phi_{1:\infty}) = \sum_{k=1}^{\infty} \pi_k f(x|\phi_k). \quad (1)$$

where $DP(\gamma, H)$ is defined as a distribution of a random probability measure G over a measurable space (Θ, \mathcal{B}) , such that for any finite measurable partition (A_1, A_2, \dots, A_r) of Θ the random vector $(G(A_1), \dots, G(A_r))$ is a Dirichlet distribution with parameters $(\gamma H(A_1), \dots, \gamma H(A_r))$.

Owing to the discrete nature and infinite dimensionality of its draws, the DP is a useful prior for Bayesian mixture models. By associating different mixture components with atoms ϕ_k , and assuming $x_i | \phi_k \stackrel{iid}{\sim} f(x_i | \phi_k)$ where $f(\cdot)$ is the kernel of the mixing components, a Dirichlet process mixture model (DPM) is obtained [18].

1) *Hierarchical DPMs*: While the DPM is suitable for the clustering of exchangeable data in a single group, many real-world problems are more appropriately modelled as comprising multiple groups of exchangeable data. In such cases it is desirable to model the observations of different groups jointly, allowing them to share their generative clusters. This so-called ‘‘sharing of statistical strength’’ emerges naturally when a hierarchical structure is implemented.

The DPM models each group of documents in a collection using an infinite number of topics. However, it is desired for multiple group-level DPMs to share their clusters. The hierarchical DP (HDP) [21] offers a solution whereby base measures of group-level DPs are drawn from a corpus-level DP. In this way the atoms of the corpus-level DP are shared across the documents; posterior inference is readily achieved using Gibbs sampling [21].

B. Modelling topic evolution over time

We now show how the described HDP based model can be applied to the analysis of temporal topic changes in a *longitudinal* data corpus.

Owing to the aforementioned assumption of a temporally locally static corpus we begin by discretizing time and dividing the corpus into epochs. Each epoch spans a certain contiguous time period and has associated with it all documents with timestamps within this period. Each epoch is then modelled separately using a HDP, with models corresponding to different epochs sharing their hyperparameters and the corpus-level base measure. Hence if n is the number of epochs, we obtain n sets of topics $\phi = \{\phi_{t_1}, \dots, \phi_{t_n}\}$ where $\phi_t = \{\theta_{1,t}, \dots, \theta_{K_t,t}\}$ is the set of topics that describe epoch t , and K_t their number.

1) *Topic relatedness*: Our goal now is to track changes in the topical structure of a data corpus over time. The simplest changes of interest include the emergence of new topics, and the disappearance of others. More subtly, we are also interested in how a specific topic changes, that is, how it evolves over time in terms of the contributions of different words it comprises. Lastly, our aim is to be able to extract and model complex structural changes of the underlying topic content which result from the interaction of topics. Specifically, topics, which can be thought of as collections of memes, can merge to form new topics or indeed split into more nuanced memetic collections. This information can provide valuable insight into the refinement of ideas and findings in the scientific community, effected by new research and accumulating evidence.

The key idea behind our tracking of simple topic evolution stems from the observation that while topics may change significantly over time, changes between successive epochs are limited. Therefore we infer the continuity of a topic in one epoch by relating it to all topics in the immediately

subsequent epoch which are sufficiently similar to it under a suitable similarity measure – we adopt the well known Bhattacharyya distance (BHD). This can be seen to lead naturally to a similarity graph representation whose nodes correspond to topics and whose edges link those topics in two epochs which are related. Formally, the weight of the directed edge that links $\phi_{j,t}$, the j -th topic in epoch t , and $\phi_{k,t+1}$ is $\rho_{\text{BHD}}(\phi_{j,t}, \phi_{k,t+1})$ where ρ_{BHD} denotes the BHD.

In constructing a similarity graph a threshold to used to eliminate automatically weak edges, retaining only the connections between sufficiently similar topics in adjacent epochs. Then the disappearance of a particular topic, the emergence of new topics, and gradual topic evolution can be determined from the structure of the graph. In particular if a node does not have any edges incident to it, the corresponding topic is taken as having emerged in the associated epoch. Similarly if no edges originate from a node, the corresponding topic is taken to vanish in the associated epoch. Lastly when exactly one edge originates from a node in one epoch and it is the only edge incident to a node in the following epoch, the topic is understood as having evolved in the sense that its memetic content may have changed.

A major challenge to the existing methods in the literature concerns the detection of topic merging and splitting. Since the connectedness of topics across epochs is based on their similarity what previous work describes as ‘splitting’ or indeed ‘merging’ does not adequately capture these phenomena. Rather, adopting the terminology from biological evolution, a more accurate description would be ‘speciation’ and ‘convergence’ respectively. The former is illustrated in Fig 1(a) whereas the latter is entirely analogous with the time arrow reversed. What the conceptual diagram shown illustrates is a slow differentiation of two topics which originate from the same ‘parent’. Actual topic splitting, which does not have a biological equivalent in evolution, and which is conceptually illustrated in Fig 1(b) cannot be inferred by measuring topic similarity. Instead, in this work we propose to employ the Kullback-Leibler divergence (KLD) for this purpose. This divergence is asymmetric can be intuitively interpreted as measuring how well one probability distribution ‘envelops’ another. KLD between two probability distributions $p(i)$ and $q(i)$ is defined as follows:

$$\rho_{\text{KLD}} = \sum_i p(i) \log \frac{p(i)}{q(i)} \quad (2)$$

It can be seen that a high penalty is incurred when $p(i)$ is significant and $q(i)$ is low. Hence, we use the BHD to track gradual topic evolution, speciation, and convergence, while the KLD (computed both in forward and backward directions) is used to detect topic splitting and merging.

2) *Automatic temporal relatedness graph construction*: Another novelty of the work first described in this paper concerns the building of the temporal relatedness graph. We achieve this almost entirely automatically, requiring only one free parameter to be set by the user. Moreover the meaning of the parameter is readily interpretable and understood by a non-expert, making our approach highly usable.

Our methodology comprises two stages. Firstly we consider all inter-topic connections present in the initial fully connected graph and extract the empirical estimate of the corresponding cumulative density function (CDF). Then we prune the graph based on the operating point on the relevant

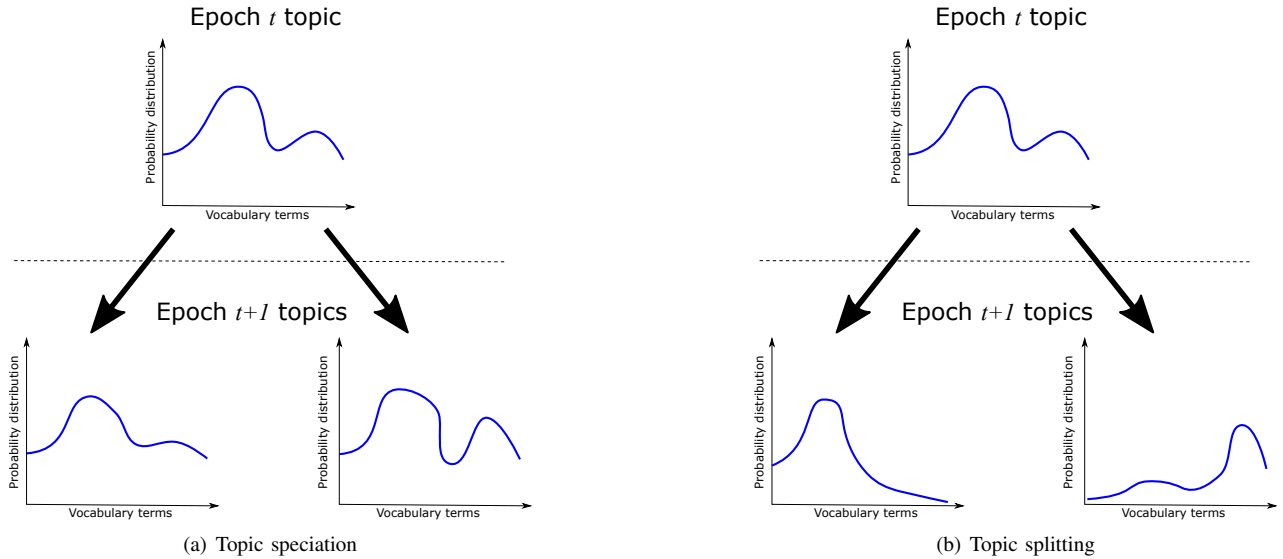


Fig. 1. This paper is the first work to describe the difference between two topic evolution phenomena: (a) topic speciation and (b) topic splitting.

CDF. In other words if F_ρ is the CDF corresponding to a specific initial, fully connected graph formed using a particular similarity measure (BHD or KLD), and $\zeta \in [0, 1]$ the CDF operating point, we prune the edge between topics $\phi_{j,t}$ and $\phi_{k,t+1}$ iff $\rho(\phi_{j,t}, \phi_{k,t+1}) < F_\rho^{-1}(\zeta)$.

III. EVALUATION AND DISCUSSION

We now analyse the performance of the proposed framework empirically on a large real world data set.

A. Evaluation data

We used the PubMed interface to access the US National Library of Medicine and retrieve from it scholarly articles. We searched for publication on the metabolic syndrome (MetS) using the keyphrase “metabolic syndrome” and collected papers written in English. The earliest publication found was that by Berardinelli *et al.* [8]. We collected all matching publications up to the final one indexed by PubMed on 10th Jan 2016, yielding a corpus of 31,706 publications.

1) *Pre-processing*: The raw data collected from PubMed is in the form of free text. To prepare it for automatic analysis a series of ‘pre-processing’ steps are required. The goal is to remove words which are largely uninformative, reduce dispersal of semantically equivalent terms, and thereafter select terms which are included in the vocabulary over which topics are learnt.

We firstly applied soft lemmatization using the WordNet[®] lexicon [17] to normalize for word inflections. No stemming was performed to avoid semantic distortion often effected by heuristic rules used by stemming algorithms. After lemmatization and the removal of so-called stop-words, we obtained approximately 3.8 million terms in the entire corpus when repetitions are counted, and 46,114 unique terms. Constructing the vocabulary for our method by selecting the most frequent terms which explain 90% of the energy in a specific corpus resulted in a vocabulary containing 2,839 terms.

B. Results

We started evaluation by examining whether the two topic relatedness measures (BHD and KLD) are capturing different

aspects of relatedness. To obtain a quantitative measure we looked at the number of inter-topic connections formed in respective graphs both when the BHD is used as well as when the KLD is applied instead. The results were normalized by the total number of connections formed between two epochs, to account for changes in the total number of topics across time. Our results are summarized in Fig 2. A significant difference between the two graphs is readily evident – across the entire timespan of the data corpus, the number of Bhattacharyya distance based connections also formed through the use of the KLD is less than 40% and in most cases less than 30%. An even greater difference is seen when the proportion of the KLD connections is examined – it is always less than 25% and most of the time less than 15%.

To get an even deeper insight into the contribution of the two relatedness measures, we examined the corresponding topic graphs before edge pruning. The plot in Fig 3 shows the variation in inter-topic edge strengths computed using the BHD and the KLD (in forward and backward directions) – the former as the x coordinate of a point corresponding to a pair of topics, and the latter as its y coordinate. The scatter of data in the plot corroborates our previous observation that the two similarity measures indeed do capture different aspects of topic behaviour.

We performed extensive qualitative analysis which is necessitated by the nature of the problem at hand and the so-called ‘semantic gap’ that underlies it. In all cases we found that our algorithm revealed meaningful and useful information, as confirmed by an expert in the area of metabolic MetS research.

Our final contribution comprises a web application which allows users to upload and analyse their data sets using the proposed framework. The application allows a range of powerful tasks to be performed quickly and in an intuitive manner. For example, the user can search for a given topic using keywords (and obtain a ranked list), trace the origin of a specific topic backwards in time, or follow its development in the forward direction, examine word clouds associated with topics, display a range of statistical analyses, or navigate the temporal relatedness graph.

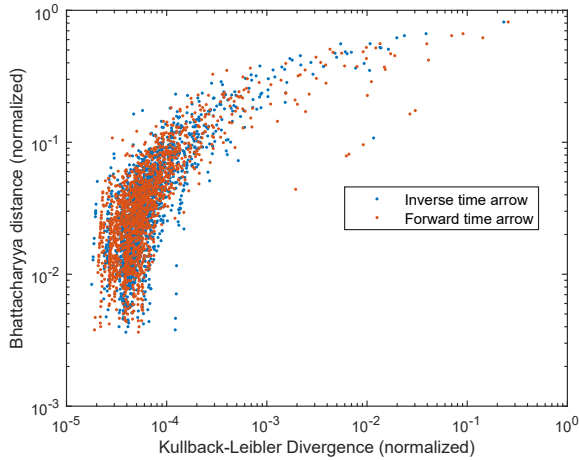
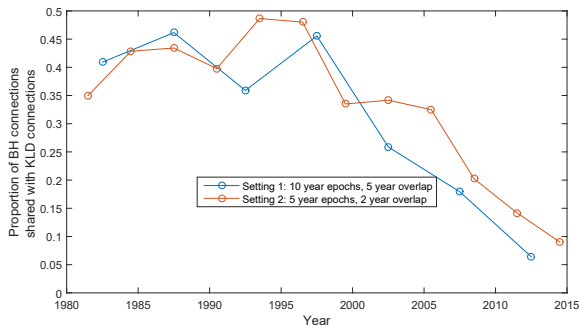
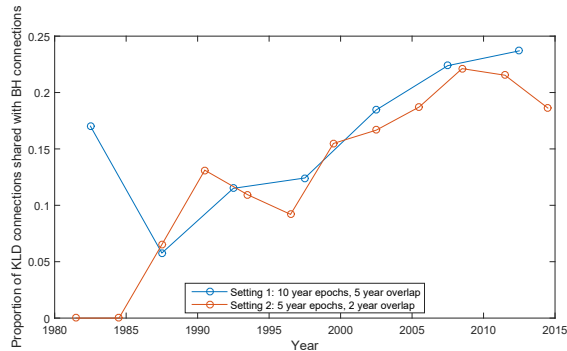


Fig. 3. Relationship between inter-topic edge strengths computed using the BHD and the KLD before the pruning of the respective graphs.



(a) BHD-KLD normalized overlap



(b) KLD-BHD normalized overlap

Fig. 2. The proportion of topic connections shared between the BHD and the KLD temporal relatedness graphs, normalized by (a) the number of BHD connections, and (b) the number of KLD connections, in an epoch.

IV. SUMMARY AND CONCLUSIONS

In this work we presented a case for the importance of use of advanced machine learning techniques in the analysis and interpretation of medical literature. We described a novel framework based on non-parametric Bayesian techniques which is able to extract and track complex, semantically meaningful changes to the topic structure of a longitudinal document corpus. Moreover this work is the first to describe and present a method for differentiating between two types

of topic structure changes, namely topic splitting and what we termed topic speciation. Experiments on a large corpus of medical literature concerned with the metabolic syndrome was used to illustrate the performance of our method. Lastly, we developed a web application which allows users such as medical researchers to upload their data sets and apply our method for their analysis; the application and its code will be made freely available following publication.

REFERENCES

- [1] O. Arandjelović. Assessing blinding in clinical trials. *Advances in Neural Information Processing Systems*, 25:530–538, 2012.
- [2] O. Arandjelović. A new framework for interpreting the outcomes of imperfectly blinded controlled clinical trials. *PLOS ONE*, 7(12):e48984, 2012.
- [3] O. Arandjelović. Clinical trial adaptation by matching evidence in complementary patient sub-groups of auxiliary blinding questionnaire responses. *PLOS ONE*, 10(7):e0131524, 2015.
- [4] O. Arandjelović. Discovering hospital admission patterns using models learnt from electronic hospital records. *Bioinformatics*, 31(24):3970–3976, 2015.
- [5] O. Arandjelović. Prediction of health outcomes using big (health) data. In *Proc. International Conference of the IEEE Engineering in Medicine and Biology Society*, pages 2543–2546, August 2015.
- [6] O. Arandjelović. Sample-targeted clinical trial adaptation. In *Proc. AAAI Conference on Artificial Intelligence*, 3:1693–1699, 2015.
- [7] O. Arandjelović. On the discovery of hospital admission patterns – a clarification. *Bioinformatics*, 2016. DOI: 10.1093/bioinformatics/btw049.
- [8] W. Berardinelli, J. G. Cordeiro, D. de Albuquerque, and A. Couceiro. A new endocrine-metabolic syndrome probably due to a global hyperfunction of the somatotrophin. *Acta Endocrinologica*, 12(1):69–80, 1953.
- [9] A. Beykikhoshk, O. Arandjelović, D. Phung, and S. Venkatesh. Hierarchical Dirichlet process for tracking complex topical structure evolution and its application to autism research literature. In *Proc. Pacific-Asia Conference on Knowledge Discovery and Data Mining*, 1:550–562, 2015.
- [10] A. Beykikhoshk, O. Arandjelović, D. Phung, and S. Venkatesh. Overcoming data scarcity of Twitter: using tweets as bootstrap with application to autism-related topic content analysis. In *Proc. IEEE/ACM International Conference on Advances in Social Network Analysis and Mining*, pages 1354–1361, August 2015.
- [11] A. Beykikhoshk, O. Arandjelović, D. Phung, S. Venkatesh, and T. Caelli. Data-mining Twitter and the autism spectrum disorder: a pilot study. In *Proc. IEEE/ACM International Conference on Advances in Social Network Analysis and Mining*, pages 349–356, 2014.
- [12] A. Beykikhoshk, O. Arandjelović, D. Phung, S. Venkatesh, and T. Caelli. Using Twitter to learn about the autism community. *Social Network Analysis and Mining*, 5(1):5–22, 2015.
- [13] D. Blei and J. Lafferty. Dynamic topic models. In *Proc. IMLS International Conference on Machine Learning*, pages 113–120, 2006.
- [14] B. Christensen and G. Ellingsen. Evaluating model-driven development for large-scale EHRs through the openEHR approach. *Int J Med Inform*, 89:43–54, 2016.
- [15] F. J. Dyson. Is science mostly driven by ideas or by tools? *Science*, 338(6113):1426–1427, 2012.
- [16] T. S. Ferguson. A Bayesian analysis of some nonparametric problems. *The Annals of Statistics*, pages 209–230, 1973.
- [17] G. A. Miller. WordNet: a lexical database for English. *Communications of the ACM*, 38(11):39–41, 1995.
- [18] M. N. Radford. Markov chain sampling methods for Dirichlet process mixture models. *Journal of Computational and Graphical Statistics*, 9(2):249–265, 2000.
- [19] RGI-CGHR Collaborators. Report on the causes of death in India: 2001–2003. *Office of the Registrar General of India*, 2009.
- [20] S. Richardson and P. J. Green. On Bayesian analysis of mixtures with an unknown number of components (with discussion). *Journal of the Royal Statistical Society*, 59(4):731–792, 1997.
- [21] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [22] I. Vasiljeva and O. Arandjelović. Prediction of future hospital admissions – what is the tradeoff between specificity and accuracy? In *Proc. International Conference on Bioinformatics and Computational Biology*, 2016.
- [23] I. Vasiljeva and O. Arandjelović. Towards sophisticated learning from EHRs: increasing prediction specificity and accuracy using clinically meaningful risk criteria. In *Proc. International Conference of the IEEE Engineering in Medicine and Biology Society*, 2016.