A NOVEL TEMPORAL TEXT MINING METHODS FOR MUSIC INFORMATION RETRIEVAL

Prabhjot Kaur¹, Er.Samandeep Singh²

¹PG Student, ²Asst. Prof., Department of CSE, Global Institute of Management & Emerging Technology, Gurdaspur (Punjab)

ABSTRACT

Detecting and tracking of temporal data is an important task in multiple applications. In this paper we study temporal text mining methods for Music Information Retrieval. We compare two Ways of detecting the temporal latent semantics of a corpus extracted from Wikipedia, using a stepwise Probabilistic Latent Semantic Analysis (PLSA) approach and a global multi-way PLSA method. The analysis indicates that the global analysis method is able to identify relevant trends which are difficult to get using a step-by-step approach. Furthermore we show that inspection of PLSA models with different number of factors may reveal the stability of temporal clusters making it possible to choose the relevant number of factors.

Keywords: Probabilistic Latent Semantic, Nonnegative Matrix Factorization, Tensor factorization.

I. INTRODUCTION

Music Information Retrieval (MIR) is a multifaceted field, which until recently mostly focused on audio analysis. The use of textual descriptions, beyond using genres, has grown in popularity with the advent of different music webs ites, e.g. "Myspace.com", where abundant data about music has become easily available. This has for instance been investigated in [1], where textual descriptions of music were retrieved from the Web to find similarity of artists. The unstructured data retrieved using web crawling produces a lot of data, which requires cleaning to produce terms that actually describe musical artists and concepts. Community based music web services such as tagging based systems, e.g. Last.fm, have also shown to be a good basis for extracting latent semantics of musical track descriptions [2]. In this study we investigate if the incorporation of time information in latent factor models enhances the detection and description of topics. Tensor methods in the context of text mining have recently received some attention using high er- order decomposition methods such as the PARAllel FACtors model [3] that can be seen as a generalization of Si ngular Value Decomposition in higher dimensional arrays. The article [3] applies tensor decomposition methods suc cessfully for topic detection in e-mail correspondence over a 12 month period. The article also employs a non-negati

vely constrained PARAFAC model forming a Nonnegative Tensor Factorization analogous to the well-known Nonnegative Matrix Factorization (NMF) [4].

Probabilistic Latent Semantic Analysis (PLSA) [5] and NMF have successfully been applied in any text analysis tas ks to find interpretable latent factors. The two methods have been shown to be equivalent [6], where PLSA has the a dvantage of providing a direct probabilistic interpretation of the latent factors. Our work therefore investigates the ex tension of PLSA to tensors.

II. TEMPORAL TOPIC DETECTION

Detecting latent factors or topics in text using NMF and PLSA has assumed an unstructured and static collection of documents. Extracting topics from a temporally changing text collection has received some attention lately, for instance by [7] and also touched by [8]. These works investigate text streams that contain documents that can be assigned a timestamp y. The timestamp may for instance be the time a news story was released, or in the case of articles describing artists it can be a time span indicating the active years of the artist. Finding the evolution of topics over time r equires assigning documents $d_1, d_2, ..., d_m$ in the collection to time intervals $y_1, y_2, ..., y_l$, as illustrated in figure 1.

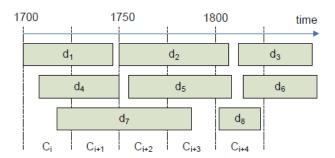


Fig.1. an example of assigning a collection of documents d_i based on the time intervals the d ocuments belong to. The assignment produces a document collection C_k for each time interval.

In contrast to the temporal topic detection approach in [7], we can assign documents to multiple time intervals, e.g. if the active years of an artist spans more than one of the chosen time intervals. The assignment of documents then provides l sub-collections C_1 , C_2 ,..., C_l of documents. The next step is to extract topics and track their evolution over time.

2.1. Stepwise temporal PLSA

The approaches to temporal topic detection presented in [7] and [8] employ latent factor methods to extract distinct t opics for each time interval, and then compare the found topics at succeeding time intervals to link the topics over time to form temporal topics.

We extract topics from each sub-collection C_k using a PLSA-model [5]. The model assumes that documents are represented as bags-of-words where each document d_i is represented by an n-dimensional vector of counts of the terms in the vocabulary; forming an nxm term by document matrix for each sub-collection C_k . PLSA is defined as a latent topic model, where documents and terms are assumed independent conditionally over topics z;

$$P(t,d)_k = \sum_{z}^{Z} P(t|z)_k P(d|z)_k P(z)_k$$
(1)

This model can be estimated using the Expectation Maximization (EM) algorithm, cf. [5]. The topic model found for each document sub-collection C_k with parameters, $\theta_k = \{P(t/z)_k, P(d/z)_k, P(z)_k\}$, need to be stringed together with the model for the next time span θ_{k+1} . The comparison of topics is done by comparing the term profiles $P(t/z)_k$ for the topics found in the PLSA model. The similarity of two profiles is naturally measured using the KL-divergence,

$$D(\theta_{k+1}||\theta_k) = \sum_{t} p(t|z)_{k+1} \log \frac{p(t|z)_{k+1}}{p(t|z)_k}$$
(2)

Determining whether a topic is continued in the next time span is quite simply chosen based on a threshold λ , such t hat two topics are linked if D $(\theta_{k+l}||\theta_k)$ is smaller than a fixed threshold λ . In this case the asymmetric KL-divergence is used in accordance with [7]. The choice of the threshold must be tuned to find the temporal links that are relevant

2.2. Multiway PLSA

The method presented above is useful to some extent, but does not fully utilize the time information that is contained in the data. Some approaches have used the temporal aspect more directly, e.g. [9] where an incrementally trainable NMF-model is used to detect topics. Using multiway models, also called tensor methods we can model the topics directly over time. The 2-way PLSA model in 1 can be extended to a 3-way model by also conditioning the topics over years *y*, as follows:

$$P(t,d,y) = \sum_{z} P(t|z)P(d|z)P(y|z)P(z)$$
(3)

The model parameters are estimated using maximum likelihood using the EM-algorithm, e.g. as in [10]. The expect ation step evaluates P(z/t, d, y) using the estimated parameters at step t. (E-step):

$$P(z|t, d, y) = \frac{p(t|z)p(d|z)p(y|z)p(z)}{\sum_{z'} p(t|z')p(d|z')p(y|z')p(z')}$$
(4)

The M-step then updates the parameter estimates.

(M-step):

$$P(z) = \frac{1}{N} \sum_{tdy} x_{tdy} P(z|t, d, y)$$
(5)

$$P(t|z) = \frac{\sum_{dy} x_{tdy} P(z|t, d, y)}{\sum_{tdy} x_{tdy} P(z|t, d, y)}$$

$$\tag{6}$$

$$P(d|z) = \frac{\sum_{ty} x_{tdy} P(z|t, d, y)}{\sum_{tdy} x_{tdy} P(z|t, d, y)}$$

$$(7)$$

$$P(y|z) = \frac{\sum_{td} x_{tdy} P(z|t, d, y)}{\sum_{tdy} x_{tdy} P(z|t, d, y)}$$
(8)

The EM algorithm is guaranteed to converge to a local maximum of the likelihood. The EM algorithm is sensitive to initial conditions, so a number of methods to stabilize the estimation have been devised, e.g. Deterministic Annealin g [5]. We have not employed these but instead rely on restarting the training procedure a number of times to find a g ood solution. The time complexity of the two PLSA approaches of course depends on the number of iterations for the method to converge. Basically the most expensive operation is the E-step of the algorithms. The cost of iteration for 2- way PLSA is O (RZ) which is calculated for each of the K time steps. In our experiments the algorithms typicall y converge in 0-50 iterations. However, the 2-way PLSA does have the advantage that the individual time steps can be calculated in parallel giving a speed-up proportional to K.

2.3. Topic model interpretation

The latent factors z of the model can be seen as topics that are present in the data. The parameters of each topic can be used as descriptions of the topic. P (t/z) represents the probabilities of the terms for the topic z, thus providing a way to find words that are representative of the topic. The most straightforward method to find these keywords is to use the words with the highest probability P (t/z). This approach unfortunately is somewhat flawed as the histogram reflects the overall frequency of words, which means that generally common words tend to dominate the P (t/z). Measuring the difference between the histograms for each topic can be measured by use of the summarized Kullback-Leible r divergence:

$$KL(z, \neg z) = \sum_{t} \underbrace{(P(t|z) - P(t|\neg z)) \log \frac{P(t|z)}{P(t|\neg z)}}_{w_t}$$
(9)

This quantity is a sum of contributions from each term t, w_t . The terms that contribute with a large value of w_t are tho se that are relatively more special for the topic z. w_t can thus be used to choose the keywords. The keywords should be chosen from the terms that have a positive value of $P(t/z) - P(t/\neg z)$ and with the largest w_t .

III. WIKIPEDIA DATA

In this experiment we investigated the description of composers in Wikipedia. This should provide us with dataset th at spans a number of years, and provides a wide range of topics. We performed the analysis on the Wikipedia data d ump saved 27th of July 2008, retrieving all documents that Wikipedians assigned to composer categories such as "B aroque composers" and "American composers". This produced a collection of 7358 documents that were parsed so t

hat only the running text was kept. Therefore words occurring in titles of documents, such as Wolfgang Amadeus M ozart, are removed from the text corpus, i.e. occurrences of the terms 'wolfgang', 'amadeus', and 'mozart' were rem oved from all documents in the corpus. Furthermore we removed irrelevant stop words based on a list of 551 words. Finally terms that occurred fewer than 3 times counted over the whole dataset and terms not occurring in at least 3 different documents were removed.

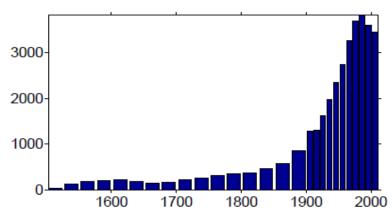


Fig. 2 Number of composer documents assigned to each of the chosen time spans.

The document collection was then represented using a bag-of-words representation forming a term-document matrix X where each element x_{td} represents the count of term t in document d. The vector x_d thus represents the term histog ram for document d. To place the documents temporally the documents were parsed to find the birth and death dates. These data are supplied in Wikipedia as documents are assigned to categories such as "1928 births" and "2007 deat hs". The dataset contains active composers from around 1500 until today. The next step was then to choose the time spans to use. We estimated the year's composers were active by removing the first 20 years of their lifetime. The resulting distribution of documents on the resulting 27 time intervals is seen in figure 2. The term by document matrix was extended with the time information by assigning the term-vector for each composer document to each era, thus forming a 3-way tensor containing terms x documents x years.

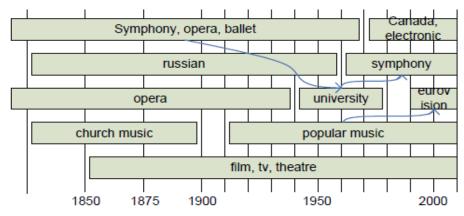


Fig. 3 Topics detected using step-by-step PLSA. The topics are depicted as connected boxes, but are the results of the KL-divergence-based linking between time slots.

The tensor was further normalized over years, such that the weight of the document summed over years is the same as in the initial term doc-matrix. I.e. $P(d) = \sum_{t,y} X_{tdy} = \sum_t X_{td}$. This was done to avoid long-lived composers dominating the resulting topics. The resulting tensor $\mathbf{X} \in \mathbb{R}^{mxnxl}$ contains 18536 terms x 7358 documents x 27 time slots with 4,038,752 non-zero entries (0.11% non-zero entries).

3.1. Term weighting

The performance of machine learning approaches in text mining often depends heavily on the preprocessing steps th at are taken. Term weighting for LSA-like methods and NMF have thus shown to be paramount in getting interpreta ble results. We applied the well-known tfidf weighting scheme, using $tf = \log(1 + x_{tdy})$ and the log-entropy document weighting, $idf = 1 + \sum_{d=1}^{D} h_{td} \log h_{td} / \log D$, where $h_{td} = \sum_{y} x_{tdy} / \sum_{dy} d_{y} x_{tdy}$. The documents in Wikipedia differ quite a lot in length, therefore we employ document normalization to avoid that long articles dominate the modeled topics.

IV. EXPERIMENTAL RESULTS

We performed experiments on the Wikipedia composer data using the stepwise temporal PLSA method and the mult iway- PLSA methods.

4.1. Stepwise temporal PLSA

The step-by-step method was trained with 5 and 16 topic PLSA models for each of the l sub-collections of document s described above. The PLSA models for each time span was trained using a stopping criterion of 10^{-5} relative chang e of the cost function, restarting the training 10 times for each model, choosing the model minimizing the likelihood. The temporal topics are produced by coupling the topics at time k and k+1 if the KL-divergence between the topic term distributions, D (θ_{k+1}/θ_k), is below a threshold λ . A low setting for λ may leave out important re relations, while a higher setting produces too many links to be interpretable. Figure 3 shows the topics found for the 20^{th} century usi ng the 5 component models. There are clearly 4 topics that are present throughout the whole period. The first is the d ifficulties in adjusting the threshold λ to find meaningful topics over time. The 5 topics that are used above do give s ome interpretable latent topics in the last decade as shown in figure 3. As an example the period 1626-1650 has the f ollowing topics:

1626-1650 41%	34%	15%	8.7%	1.4%
keyboard	madrigal	viol	baroque italy poppea italian lincoronazion opera finta era venice teatro	anglican
organ	baroque	consort		liturgi
surviv	motet	lute		prayer
italy	continuo	england		respons
church	monodi	charles		durham
nuremberg	renaissance	royalist		english
choral	venetian	masqu		chiefli
baroque	style	fretwork		england
germani	cappella	charles's		church
collect	itali	court		choral

The topics found here are quite meaningful in describing the baroque period, as the first topic describes the church music, and the second seems to find the musical styles, such as madrigals and motets. The last topic on the other han d only has a topic weight of P(z) = 1.4%. This tendency was even more distinct when using 16 components in each t ime span.

4.2. Multi-way PLSA

We then model the full tensor, described in section 3, using the mw-PLSA model. Analogously with the stepwise te mporal PLSA model we stopped training reaching a change of less than 10^{-5} of the cost function. The time evolution of topics can be visualized using the parameter P (y/z) that gives the weight of a component for each time span. Figu re 4 shows the result for 4, 8, 16, 32 and 64 components as a "heatmap", where darker colors correspond to higher v alues of P(y/z). The skewed distribution of documents over time which we described earlier emerges clearly in the di stribution of topics, as most of the topics are placed in the last century.

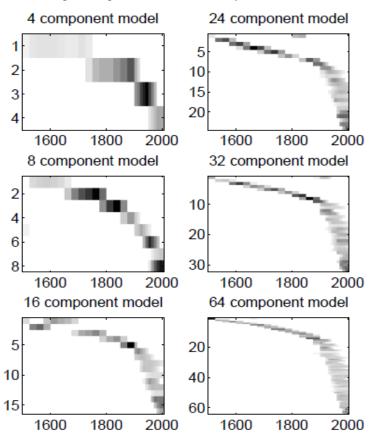


Fig. 4 Time view of components extracted using mw-PLSA, showing the time profiles P (y/z) as a he atmap. A dark color corresponds to higher values of P (y/z).

The keywords extracted using the method mentioned above are shown in table 1, showing 5 of the topics extracted by the 32 component model, including the time spans that they belong to. The first topic shown in table 1 is one of the two topics that accounts for the years 1626-1650, the keywords summarize the five topics found using mw-PLSA. T

he second topic has the keywords ragtime and rag, placed in the years 1876-1921, which aligns remarkably well with the genre description on Wikipedia: "Ragtime [...] is an originally American musical genre which enjoyed its peak popularity between 1897 and 1918". The next topic seems to describe World War II, but also contains the neoclassical movement in classical music.

Table-1 Keywords for 5 of 32 components in an mw-PLSA model. The assignment of years is given from P(y/z) and percentages placed at each column are the corresponding component weights, P(z)

1601-1700	1676-1776	1876-1921	1921-1981	1971-2008
2.10%	2.40%	4.70%	4.80%	6.40%
baroque	baroque	ragtime	concerto	single
italian	opera	sheet	nazi	chart
church	sonata	rag	war	album
continuo	italian	weltemignon	symphony	release
survive	court	nunc	piano	hit
court	harpsichord	ysa	ballet	track
organ	italy	schottisch	neoclassical	sold
motet	violinist	dimitti	neoclassic	demo
madrigal	church	blanch	choir	fan
cathedral	organ	parri	hochschul	pop

The 16 component stepwise temporal PLSA approach finds a number of topics from 1921-1940 that describe the war, such as a topic in 1921-1930 with keywords: war, time, year, life, influence and two topics in 1931-1941, 1. time, war, year, life, style and 2: theresienstadt, camp, auschwitz, deport, concentration, nazi. These are quite unrelated to music, so it is evident that the global view of topics employed in the mw-PLSA model identifies neoclassicism to be the important keywords compared to topics from other time spans. Some of the topics do overlap in time, such as the first two presented in table 1, and it is clear that they present different aspects of the music in the Baroque era, one representing church music (organ and madrigals), while the other describes opera and sonatas.

V. MULTI-RESOLUTION TOPICS

The use of different number of components in the mw-PLSA model, as seen in figure 4, shows that the addition of to pics to the model shrinks the number of years they span.

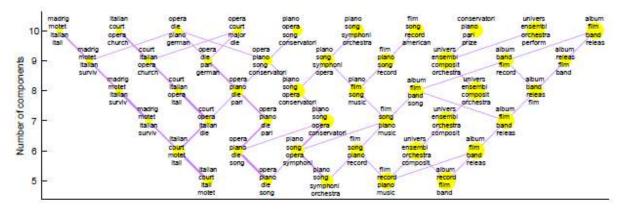


Fig. 5 The Cluster bush visualization of the results of the mw-PLSA clustering of composers.

The size of the circles correspond to the weight of the cluster, and the thickness of the line between circles how relat ed the clusters are. Each cluster is represented by the keywords and is placed according to time from left to right.

The higher specificity of the topics when using more components gives a possibility to "zoom" in on interesting topi c, while the low complexity models can provide the long lines in the data. The result for the mw-PLSA-based cluster ing is shown in figure 5. The clusters are sorted such that the clusters placed earliest in time are placed left. The clust er bush could therefore be good tool for exploring the topics at different time spans to get an estimate of the number of components needed to describe the data.

VI. CONCLUSION AND FUTURE SCOPE

We have investigated the use of time information in analysis of musical text in Wikipedia. It was shown that the use of time information produces meaningful latent topics, which are not readily extractable from this text collection wit hout any prior knowledge. The stepwise temporal PLSA approach is quite fast to train and processing for each time span can readily be processed in parallel. The multiway PLSA was shown to provide a more flexible and compact re presentation of the temporal data than stepwise temporal PLSA method. The global model would also make it possib le to do model selection overall time steps directly. The use of Wikipedia data also seems to be a very useful resource for semi-structured data for Music Information Retrieval that could be investigated further to harness the full potential of the data.

REFERENCES

- [1] P. Knees, E. Pampalk, and G. Widmer, "Artist classification with web-based data," in Proceedings of ISMIR, Ba rcelona, Spain, 2004.
- [2] M. Levy and M. Sandler, "Learning latent semantic models for music from social tags," *Journal of New Music R esearch*, vol. 37, No. 2, pp. 137–150, 2008.
- [3] B. Bader, M. Berry, and M. Browne, *Survey of Text Mining II Clustering, Classification, and Retrieval,* chapter Discussion tracking in Enron email using PARAFAC, Springer, 2008.
- [4] D. D. Lee and H. S. Seung, "Learning the parts of objects by non-negative matrix factorization," *Nature*, vol. 40 1, no. 6755, pp. 788–791, 1999.
- [5] T. Hofmann, "Probabilistic Latent Semantic Indexing," in Proc. 22nd Annual ACM Conf. on Research and Devel opment in Information Retrieval, Berkeley, California, August 1999.
- [6] E. Gaussier and C. Goutte, "Relation between plsa and nmf and implications," in Proc. 28th annual ACM SIGIR conference, New York, USA, 2005.
- [7] Q. Mei and C. Zhai, "Discovering evolutionary theme patterns from text: an exploration of Temporal text mining," *in Proc. of KDD* 05. 2005.
- [8] M. W. Berry and M. Brown, "Email surveillance using nonnegative matrix factorization," *Computational and M athematical Organization Theory*, vol. 11, pp. 249–264, 2005.

- [9] B. Cao, D. Shen, J-T Sun, X. Wang, Q. Yang, and Z. Chen, "Detect and track latent factors with online nonnegat ive matrix factorization," *in Proc. of IJCAI-07*, Hyderabad, India, 2007.
- [10] K. Yoshii, M. Goto, K. Komatani, T. Ogata, and H. G. Okuno, "An efficient hybrid music recommender system using an incrementally trainable probabilistic generative model," *Audio, Speech, and Language Processing, IE EE Transactions on*, vol. 16, no. 2, Feb. 2008.
- [11] F. A. Nielsen, D. Balslev, and L. K. Hansen, "Mining the posterior cingulate: Segregation between memory and pain components," *NeuroImage*, vol. 27, no. 3, pp. 520–532, 2005.
- [12] Jiban K. Pal, "Usefulness and applications of data mining in extracting information from different perspectives" Annals of library and information studies, Vol-58, March 2011.
- [13] M. Banu Priya, Dr.A.Kumaravel, "Methodologies for Trend Detection Based on Temporal Text Mining", *IJCS MC*, Vol. 2, Issue. 4, pg.540 554, April 2013.