

SOCIAL MEDIA TO MONITOR PEOPLE'S HEALTH OVER TIME

A.Soumya¹, B.Sharmila²

^{1,2}Assistant Professor

^{1,2}Department of Computer Science & Engineering

^{1,2}Vignan Institute of Technology & Science, Hyderabad

ABSTRACT:

Social media has become a major source for the study of every aspect of everyday life. Public health can now be analyzed on Twitter due to committed latent thematic research tools such as the Ailment Subject Aspect Model (ATAM). We're interested in using social media in this research to track people's health over time. The use of tweets has several benefits including instantaneous data availability at virtually no cost. Early monitoring of health data is complementary to post-factum studies and enables a range of applications such as measuring behavioral risk factors and triggering health campaigns. We formulate two problems: health transition detection and health transition prediction. We first propose the Temporal Ailment Topic Aspect Model (TM-ATAM), a new latent model dedicated to solving the first problem by capturing transitions that involve health-related topics. TM-ATAM is a non-obvious extension to ATAM that was designed to extract health-related topics. It learns health-related topic transitions by minimizing the prediction error on topic distributions between consecutive posts at different time and geographic granularities. To solve the second problem, we develop T-ATAM, a Temporal Ailment Topic Aspect Model where time is treated as a random variable natively inside ATAM. Our experiments on an 8-month corpus of tweets show that TM-ATAM outperforms TM-LDA in estimating health-related transitions from tweets for different geographic populations. We are looking at TM-ATAM's ability to distinguish shifts in different geographic regions due to climate conditions. We then explain how to use T-ATAM to forecast the most significant change and equate T-ATAM with CDC data and Google Flu patterns.

1.INTRODUCTION

Social media has become an important source of information for the study of many facets of everyday life. In particular, monitoring of public health can be carried out on Twitter to assess the well-being of the different geographical populations. The ability to model transitions for ailments and detect statements such as "people talk about smoking and cigarettes before talking about respiratory problems", or "people talk about headaches and stomach ache in any order", has a range of applications in syndromic surveillance such as measuring behavioral risk factors and triggering public health campaigns. Popular probabilistic topic modeling methods such as Latent Dirichlet Allocation [2] and pLSA [4] have a long history of successful application to news articles and academic abstracts. However, the small size of social media content poses serious challenges to the efficacy of such methods [9]. Dedicated methods, such as the Ailment Topic Aspect Model (ATAM), have thus been proposed to discover ailments from tweets [5]. While the primary goal of probabilistic topic modeling is to learn topic models, an equally interesting objective is to examine topic transitions. A temporal extension to LDA (TM-LDA) was hence developed for discovering the evolution of general-purpose topics in tweets [8]. In this paper, we examine the feasibility of measuring and predicting ailment transitions in Twitter, by combining ATAM and TM-LDA into a new model, coined TM-ATAM. Our model is different from dynamic topic models such as [1,7], as it is designed to learn topic transition patterns from temporally-ordered posts, while dynamic topic models focus on changing word distributions of topics over time. TM-ATAM learns transition parameters by minimizing the prediction error on ailment distributions of consecutive periods at different temporal and geographic granularities. The effectiveness of TM-ATAM requires to carefully model two key granularities, temporal and geographic. A temporal granularity that is too-fine

may result in sparse and spurious transitions whereas a too-coarse one could miss valuable ailment transitions. Similarly, a too-fine geographic granularity may produce false positives and a too coarse one may cover a user population that is exposed to different weather conditions and miss meaningful transitions. Our experiments on a corpus of more than 500K health-related and geo-localized tweets collected over a period of 8 months, show that TM-ATAM outperforms ATAM, TM-LDA and LDA in estimating temporal health-related topic transitions of different geographic populations. The health-related topic transitions we unveiled can be broadly classified in 2 kinds: stable-topics are those where a health-related topic is mentioned continuously. One-way-transitions cover the case where some topics are discussed after others. For example, our study of tweets from Arizona revealed many self-transitions such as headaches and body pain. On the other hand, tweets about smoking, drugs and cigarettes in California, are followed by respiratory ailments.

II.MODEL, PROBLEM AND APPROACH

Table 1 summarizes the terminology we use throughout this paper. By using suitable geographic granularity g (country, state, county) and temporal granularity t (week, biweek and months), we build our document sets $D_{t g}$. While LDA is successful at uncovering generic topics, its limitations at discovering infrequent and specific topics such as health has already been shown [5]. The probabilistic Ailment Topic Aspect Model (ATAM) was designed specifically to uncover latent health-related topics present in a collection of tweets [5]. ATAM achieves remarkable improvement over LDA in discovering topics that correspond to ailments (in addition to discovering general topics). The topic distribution vector generated by ATAM for a sample tweet is shown in Figure 1. Note the stronger relevance to health-related matters in this vector than in the topic distribution vector generated by LDA for the same tweet. While ATAM is effective at modeling health-related topics, it is not designed to model topic transitions over time.

Table 1: Mapping tweets to documents

Term	Description
\mathcal{P}	set of (tweet) posts
\mathcal{G}	set of regions
\mathcal{T}	set of time periods
\mathcal{P}_g^t	posts from region g during time t
D_g^t	document-set built by mapping the content of each post $p \in \mathcal{P}_g^t$ to a document
Θ_g^t	ailment distribution vector for document-set D_g^t of region g during time t
m	distance measure between distributions

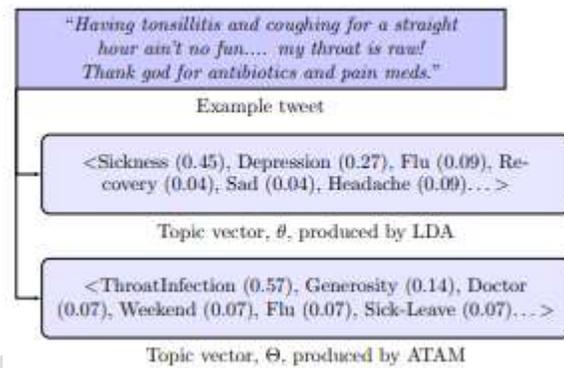


Figure 1: LDA vs ATAM: Comparison of topic distributions for an example tweet

2.1 Ailment prediction problem

In [8], TM-LDA was introduced to extend LDA with modeling topic evolution over time. However, While being quite elegant in modeling general-purpose topics TM-LDA is not specialized to capture health transitions over time. Let $\Theta_{t g}$ be a ailment distribution vector where the weight of each ailment is representative of the discourse density of ailment in the tweets originating from region g during period t . For a region g , the interval of time spanning a set of consecutive time periods $\{t_i, t_{i+1}, \dots\}$ during which discovered ailment distributions $\{\Theta_{t_i g}, \Theta_{t_{i+1} g}, \dots\}$ do not change appreciably forms a homogenous time period w.r.t. ailments. By definition, a homogenous time period is (nearly) homogeneous in terms of ailments. In other words, the ailments evolve in a smooth fashion within a homogenous time period and change abruptly across homogenous time period boundary. We posit that such homogenous time periods exist after which they encounter change-points

in ailment topic discussions. These change-points in ailment topic discussions may be caused by onset of the disease or some other external factors. Nevertheless, they are the interesting points for analyzing purposes. As an example, in Figure 2, we show the difference between ailment distributions of consecutive months for 3 different regions Kuala Lumpur (a city in Indonesia), Oklahoma (a state in the USA), and Bristol (a city in the UK). The sharp peaks obtained validate the existence of time intervals that are homogeneous w.r.t. ailments.

Algorithm 1 TM-ATAM: *change-point* Detection and Training Ailment Distribution Predictor

```

1: for all  $g \in G$  do
2:   Run ATAM on  $D_g$ 
3:   for all  $t \in \mathcal{T}$  do:
4:     for all  $z \in \mathcal{Z}$  do:
5:        $\Theta_g^t[z] \leftarrow 0$ 
6:     end for
7:     for all  $d \in D_g^t$  do:
8:       for all  $w \in d$  do:
9:          $z \leftarrow \text{topic}(w)$ 
10:         $\Theta_g^t[z] \leftarrow \Theta_g^t[z] + \frac{1}{|d| \times |D_g^t|}$ 
11:      end for
12:    end for
13:  end for
14:   $t_c = \text{argmax}_t m(\Theta_g^{t-1}, \Theta_g^t)$ 
15:   $pre = [t_1, t_{c-1}]$ 
16:   $post = [t_c, t_{|\mathcal{T}|}]$ 
17:  for all  $s \in \{pre, post\}$  do:
18:     $A_g^s \approx A_g^{s-1} . M$ 
19:     $M = (A_g^{s-1 \top} A_g^{s-1})^{-1} A_g^{s-1 \top} A_g^s$ 
20:  end for
21: end for

```

Our problem: Given a set of documents D_{t-1}^g formed by tweets originating from a region $g \in G$ during time period $t-1$, predict Θ_{t-1}^g , the ailment distribution of documents in D_{t-1}^g , corresponding to posts from g in period $t-1$ from Θ_{t-1}^g , the ailment distribution of document D_{t-1}^g corresponding to posts from g during period $t-1$.

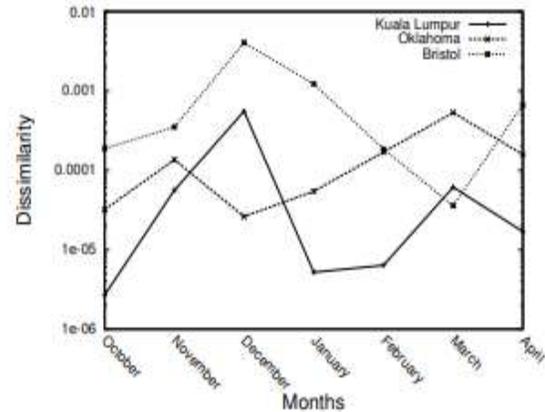


Figure 2: Topic transitions over time.

2.2 Modeling Health Topics over Time with TM-ATAM

To solve our problem, we propose TM-ATAM that builds on top of ATAM and TM-LDA. We first convert inferences of ATAM over a single document to associate with a given set of documents D_{t-1}^g , an ailment distribution, Θ_{t-1}^g . We then go on to find homogenous time periods. We model ailment transitions within each homogenous time period and when a change-point is encountered we update these transitions. This is a fresh departure from existing solutions that operate in a homogenous time period-agnostic fashion [8]. TM-ATAM, at its heart, solves the following equation.

$$A_g^t \approx A_g^{t-1} . M^* \quad (1)$$

where

$$A_g^{t-1} = \begin{pmatrix} \Theta_g^1 \\ \vdots \\ \Theta_g^t \end{pmatrix}, A_g^t = \begin{pmatrix} \Theta_g^2 \\ \vdots \\ \Theta_g^{t+1} \end{pmatrix} \quad (2)$$

M^* is an unknown transition matrix which is obtained by solving the following least squares problem.

$$M^* = \underset{M}{\text{argmin}} \|A_g^t - A_g^{t-1} . M\|_F$$

Algorithm 1 contains the steps of our solution. It has two parts: change-point detection and ailment prediction.

Change Point Detection.

We use Z to refer to the set of all health-related and nonhealth related topics. For each region $g \in G$ (Line 1) we first run ATAM over the full time period Dg (Line 2). Next for each period $t \in T$ (Line 3), we use the output of ATAM over Dg to generate a topic distribution $\Theta t g$ (Lines 4– 12). We then examine the Bhattacharya Distance between consecutive distributions $\Theta t-1 g$ and $\Theta t g$ of the region g to identify the most significant change-point , t_c , for region g (Line 14). The time periods preceding and succeeding change-point are termed as homogenous time periods .

Ailment Prediction.

In the second module of TM–ATAM algorithm, we predict distribution of ailments in twitter discourse ahead of time for each homogenous time period . Lines 17–20 of Algorithm 1 outline the steps undertaken to identify the detection of ailments for intra-homogeneous periods.

III.EXPERIMENTS

We conducted experiments to evaluate the performance of TM–ATAM and to compare with the state of the art.

3.1 Experimental setup

We employ Twitter’s Streaming API to collect tweets between 2014-Oct-8 and 2015-May-31. Collected tweets were subjected to two pre-processing steps as follows. Identifying health-related tweets: We filter the tweets returned by the Decahose Stream to obtain health-related tweets. We say that a tweet is health-related if it contains a health keyword and passes our classification criteria. The process is automated with the help of an SVM classifier [3]. To this end, 5128 tweets were annotated through crowdsourcing efforts. The precision and recall of the classifier are 0.85 and 0.44. Table 2 shows that out of the 1.36B tweets we collected, 698K were health-related. Identifying geolocalized tweets: The ability to operate seamlessly at varying geographic

resolutions mandates that the exact location of each tweet be known to TM–ATAM. Twitter affords its users the option to share their geolocation. In our case, over half a million tweets are retained (569K as indicated in Table 2). We examine various choices for the geographic granularities, temporal granularities and distance measures. TM–ATAM performs better on smaller regions and finer temporal granularity. We attribute this result to the fact that tweets from smaller regions and finer temporal granularity show less diversity in topics. We also used 2 distance measures to measure distribution difference namely, cosine similarity and bhattacharya distance. We observed that number of tweets at a given time granularity t may affect the performance of Cosine Similarity. Finally, we chose to work with geographic granularity of states , temporal granularity of months and distance measure of bhattacharya . Test-bench and measures: We run our experiments on a 32 core Intel Xeon @ 2.6Ghz CPU (with 20MB cache per core) system with 128 Gig RAM running Debian GNU/Linux 7.9 (wheezy) operating system. All subsequently discussed components were implemented in Java 1.8.0_60. We used perplexity to compare between models [6].

IV.EXPERIMENTAL RESULTS

Recall that the terms change-point and homogenous time period refer to the point in time at which discourse density of ailments changes substantially, and the time period before and after that point, respectively. We divide each homogenous time period into training and test sets. ATAM is then re-run over the training set of each homogenous time period . We then model a transition matrix M^* t_{matam} on the training set of each homogenous time period as described in Section 2.2. We compute the probability of "health topic" z for each tweet p of the first month ($|T| - 1$) in the test set using the following formulas:

$$P(z|t_{|T|-1}) = \frac{\sum_{p \in t_{|T|-1}} P(z|p \text{ for } t_{|T|-1})}{\#tweets \text{ for } t_{|T|-1}} \quad (3)$$

$$P(z|p) = \sum_{w} P(z|w)P(w|p) = \sum_{w} \frac{n(z,w)}{n(w)} P(w|p) \quad (4)$$

Here, values for $n(z, w), n(w)$ are taken from ATAM run on the training months. $P(w|p)$ is simply the number of times word w occurs in the tweet p divided by the total number of words in p . We then predict the future probability of each topic in the second month of the test data using the corresponding transition matrix $M^* \text{tmam}$. Probability of word $p_l(w_i)$ for any document set is calculated as follows:

$$p_l(w_i) = \sum_z P(w|z)P(z) = \sum_z \frac{n(z, w)}{n(z)} P(z) \quad (5)$$

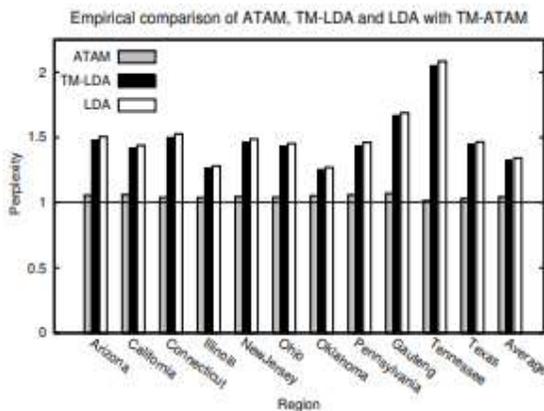


Figure 3: Comparison for top 10 active regions. Histograms denote ratio of perplexities. TM-ATAM is always at 1.0.

Having computed $P(w)$, we can compute perplexity against the words of the tweets of second test month.

TM-ATAM vs ATAM, TM-LDA vs LDA

Figure 3 shows the perplexity ratio of TM-ATAM with state-of-the-art models. If ratio computes to be less than "1" for competitor topic model, TM-ATAM is performing worse. If ratio is more than "1" for competitor topic model, TM-ATAM is performing better. In order to assert the fact that health topics transit from one to another, we compute the perplexity of ATAM on words of the first month of the test set and not predicting any topic distribution using any transition matrix. Hence, this denotes the case where health topics stay static. As shown in Figure 3, TM-ATAM beats ATAM in all social media active regions (an active region is a region where the proportion of tweets is high enough). For training TM-LDA, we merge the

training data (same as TM-ATAM) of each homogenous time period in each region and train a transition matrix of TM-LDA by solving least squares problem. For each tweet p of the first month of the test set, we compute the probability of topics using LDA trained on merged training data (Formula 3). We then predict the future probability of each topic in the following month using $M^* \text{tmlda}$. We can then compute the perplexity of TM-LDA against words of actual tweets in the test set. Figure 3 shows that TM-ATAM consistently beats TM-LDA and LDA in predicting future health topics on the test month. Perplexity is indeed lower for all words of the test month in all active states.

Homogenous Time Periods In Figure 4 we show the top-2 sharpest change-points for the top regions. Those points can be explained with weather changes in those regions. Texas can be explained with a drop in temperature while Jervis Bay can be explained by an increase in rainfall. Dublin sees its lowest temperature in November. Singapore and Manila have very similar weather conditions and exhibit the same change point.

Topic Transitions Entry m_{ij} in the transition matrix M produced by TM-ATAM, shows the degree that topic z_i will contribute to topic z_j in the subsequent homogenous time period. We adapt the threshold used in [8] to our settings: $\text{Threshold} = \mu + 2 \times \sigma_{\text{non-diagonal}}$. We identify two kinds of interesting transitions based on the above threshold: selftransitions : popular topics and one way transitions : i th topic is discussed before j th topic. As an example, one-way-transitions of California are summarized in Table 3

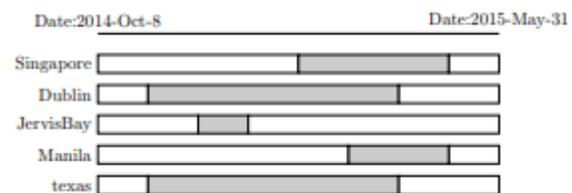


Figure 4: Top-2 Monthly homogenous time period for top active regions.

Table 3: One-Way Transitions for California (threshold: 0.815)

From Topic	To Topic	Weight
smoking/junkies /drugs/cigarettes	respiratory diseases	2.70
depression/complaining /cursing/slangs/self-pity	joint pains/body pains	3.25

Media Using Topic Models. In ECIR, pages 338–349, 2011.

V.CONCLUSION

We learned how to detect ailment distributions in social media over time. We have proposed a granularity-based model for conducting field-specific research leading to the identification of time intervals characterizing homogeneous discourse of ailments, by region. We modeled the evolution of diseases within each homogeneous area and tried to predict ailments. The fine-grained complexity of our model results in major improvements over state-of - the-art methods.

REFERENCES

- [1] D. M. Blei and J. D. Lafferty. Dynamic Topic Models. In ICML, pages 113–120, 2006.
- [2] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet Allocation. JMLR, 3:993–1022, 2003.
- [3] C. Cortes and V. Vapnik. Support-vector networks. Machine Learning, 20(3):273–297, 1995.
- [4] T. Hofmann. Probabilistic Latent Semantic Indexing. In SIGIR, pages 50–57, 1999.
- [5] M. J. Paul and M. Dredze. You Are What You Tweet: Analyzing Twitter for Public Health. In ICWSM, 2011.
- [6] H. M. Wallach, I. Murray, R. Salakhutdinov, and D. Mimno. Evaluation methods for topic models. In ICML, pages 1105–1112, 2009.
- [7] X. Wang and A. McCallum. Topics Over Time: A Non-Markov Continuous-time Model of Topical Trends. In KDD, pages 424–433, 2006.
- [8] Y. Wang, E. Agichtein, and M. Benzi. TM-LDA: Efficient Online Modeling of Latent Topic Transitions in Social Media. In KDD, pages 123–131, 2012.
- [9] W. X. Zhao, J. Jiang, J. Weng, J. He, E. Lim, H. Yan, and X. Li. Comparing Twitter and Traditional