

# Detecting and Characterizing Events

Allison J.B. Chaney  
Princeton University  
achaney@cs.princeton.edu

Hanna Wallach  
Microsoft Research  
wallach@microsoft.com

David M. Blei  
Columbia University  
david.blei@columbia.edu

Matthew Connelly  
Columbia University  
mjc96@columbia.edu

## ABSTRACT

Significant events are characterized by interactions between entities (e.g., countries, organizations, individuals) that deviate from typical interaction patterns. Investigators, such as historians, commonly read large quantities of text to construct an accurate picture of who, what, when, and where an event happened. In this work, we present the *Capsule* model for analyzing documents to identify and characterize events of potential significance. Specifically, we develop a model based on topic modeling to distinguish between topics that describe “business-as-usual” and topics that deviate from these patterns. To demonstrate this model, we analyze real-world datasets, including a corpus of over 2 million US State Department cables from the 1970s; we provide open-source implementations of an inference algorithm for the Capsule model and a visualization of its results.

## CCS Concepts

•Information systems → Document collection models; •Human-centered computing → Visualization toolkits;

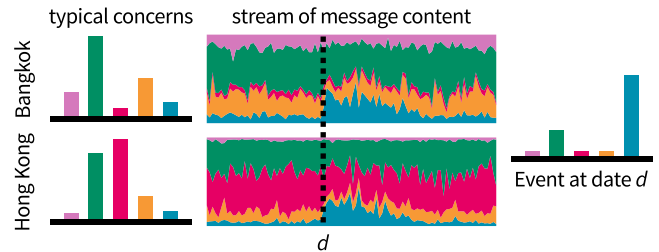
## Keywords

ACM proceedings; L<sup>A</sup>T<sub>E</sub>X; text tagging

## 1. INTRODUCTION

Historical events are difficult to define; historians and political scientists read large quantities of text to construct an accurate picture of a single event. Events are interesting by definition: they are the hidden causes of anomalous observations. But they are also inherently abstract—we can observe that changes occur, but we cannot directly observe whether or not an event occurs.

Consider embassies sending diplomatic messages, such as shown in Figure 1. The Bangkok and Hong Kong embassies have *typical concerns* about which they usually send messages. At date  $d$ , however, the message content changes for



**Figure 1: Cartoon intuition of Capsule.** Both the Bangkok and Hong Kong embassies have typical concerns about which they usually send messages (represented in topic space). When an event occurs at date  $d$ , the stream of message content alters to include the event, then fades back to “business as usual.” Capsule discovers both entities’ typical concerns and the event locations and content.

both embassies—again, we only observe the changes in message content, and do not observe the event directly. Our first goal is to determine *when* events happen, or identify these rare but pervasive deviations from the typical concerns.

Our second goal is to characterize *what* occurs. We rely on topic models [3] to summarize documents and use that same latent space to characterize events.

We develop a Bayesian model that discovers the typical concerns of authors, identifies when events occur, and characterizes these events; we call this the *Capsule* model, as it encapsulates events.

Our final goal is to visualize the results of the Capsule model to make them accessible. We provide source code for both Capsule and its associated visualization.

We first review previous research related to event detect, summarization, and visualization. In Section 2, we describe the Capsule model and how to infer the latent parameters (the appendix provides further inference details). Section 3 provides an exploration of results on simulated and three real-world datasets, and we conclude with a discussion in Section 4.

**Related work.** We first review previous work on automatic event detection and other related concepts.

While Capsule uses text documents and associated metadata as input, event detection is often performed with univariate input data. In this context, bursts that deviate from typical behavior (e.g., noisy constant or a repeating pattern) can define an event [21, 18]; Poisson Processes [20] are often used to model events under this definition. Alternatively,

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

KDD '16 San Francisco, California USA

© 2016 Copyright held by the owner/author(s).

ACM ISBN ??...\$15.00

DOI: ?

events can be construed as “change points” to mark when typical observations shift semi-permanently from one value to another [16]. In both univariate and multivariate settings, the goal is often the same: analysts want to predict whether or not a rare events will occur [40, 9]. Capsule, in contrast, is designed to help analysts explore and understand the original data: our goal is interpretability, not prediction.

Text is often used in event detection, as it is an abundant source of data. In some applications, documents themselves are considered to be observed events [28, 30], or events are predetermined and tracked through the documents [41, 37]. We are interested in detecting *unobserved* events which can be characterized by patterns in the data.

A common goal is to identify clusters of documents; these approaches are used on news articles [44, 43, 42, 24, 39, 1] and social media posts [37, 23, 19, 34, 32, 2, 35]. In the case of news articles, the task is to create new clusters as novel news stories appear—this does not help disentangle typical content from rare events of interest. Social media approaches identify rare events, but the methods are designed for short, noisy documents; they are not appropriate for larger documents that contain information about a variety of subjects.

Many existing methods use document terms as features, frequently weighted by tf-idf value [13, 22, 5, 10, 43, 44]; here, events are bursts in groups of terms. Because language is high dimensional, using terms as features limits scalability.

Topic models [3] reduce the dimensionality of text data; they have been used to help detect events mentioned in social media posts [23, 12] and posts relevant to monitored events [37]. We rely on topic models to characterize both typical content and events, but grouped observations can also be summarized directly [30, 7, 14].

In addition to text data over time, author [43], news outlet [39], and spatial information [29, 27, 26] can be used to augment event detection. Capsule uses author information in order to characterize typical concerns of authors.

Detecting and characterizing relationships [36, 25, 10] is related to event detection. When a message recipient is known, Capsule’s author input can be replaced with a sender-receiver pair, but the model could be further tailored for interactions within networks.

Once events have been identified and characterized, visualization translates a model’s output into sometime interpretable for non experts. LeadLine [12] is an excellent example of a visualization of event detection. We build on topic model visualization concepts [8] to provide tailored visualization code for Capsule.

## 2. THE CAPSULE MODEL

In this section we develop the Capsule model. Capsule captures patterns in entity behavior and identifies events that cause deviations from these patterns among many entities. The model relies on rich entity behavior data over time, such as messages being sent between entities; text data can be summarized (making the model more tractable) with a topic model [3]. We first review topic models at a high level and give the intuition on Capsule. Then, we formally specify our model and discuss how we learn the hidden variables.

**Background: Topic Models.** Capsule relies on topic models to summarize text data, making the model tractable. Topic models are algorithms for discovering the main themes in a large collection of documents; each document can then

be summarized in terms of the global themes. More formally, a topic  $k$  is a probability distribution over the set of vocabulary words. Each document  $d$  is represented as a distribution over topics  $\theta_d$ . Thus we can imagine that when we generate a document, we first pick which topics are relevant (and in what proportions); then, for each word, we select a single topic from this distribution over topics, and finally select a vocabulary term from the corresponding topic’s distribution over the vocabulary. We use the LDA topic model [4, 17] to summarize text data, and assume that these summaries are held fixed. Our model could be extended to include topic modeling as component, but in practice the results would be similar to the stage-wise approach.

**The Capsule Model.** Topic models are often applied to provide a structure for an otherwise unstructured collection of documents. Documents, however, are often accompanied by metadata, such as the date written or author attribution; this information is not exploited by traditional topic models. The Capsule model uses both author and date information to identify and characterize events that influence the content of the collection.

Consider an entity like the Bangkok American embassy, shown in Figure 1. We can imagine that there is a stream of messages (or *diplomatic cables*) being sent by this embassy—some might be sent to the US State Department, others to another American embassy like Hong Kong. An entity will usually talk about certain topics; the Bangkok embassy, for instance, is concerned with topics regarding southeast Asia more generally.

Now imagine that at a particular time, an event occurs, such as the capture of Saigon during the Vietnam war. We do not directly observe that events occur, but each event can again be described in the same topic space used to describe individual messages. Further, when an event occurs, the message content changes for multiple entities. The day following the capture of Saigon, the majority of the diplomatic cables sent by the Bangkok embassy were about Vietnam war refugees. Thus we imagine that an entity’s stream of messages is controlled by what it usually talks about as well as the higher level stream of unobserved events.

**Model Specification.** We formally describe Capsule. The observed data are documents represented in topic space; each document has an author (or entity) and a time (or date) associated with it. The document content for each document  $d$  is represented as  $\theta_d$ , a  $K$ -dimensional vector, where  $K$  is the number of topics. The author and time associated with document  $d$  are represented as  $n_d$  and  $m_d$ , respectively.

The hidden variables of this model are the authors’ typical concerns, event occurrences, and event descriptions. We represent the concerns of author  $n$  as  $\phi_n$ , also a  $K$ -dimensional topic vector. For each time  $t$  we represent whether or not an event occurs with  $\epsilon_t$ ; when an event does occur we represent its content as  $\pi_t$ , another  $K$ -dimensional topic vector.

Conditional on the hidden variables and the author and time metadata, Capsule is a model of how document topics  $\theta_d$  came to be; we generate the topics for each document

$$\theta_{d,k} \sim \text{Gamma} \left( \phi_{n_d,k} + \sum_{t=1}^T f(t, m_d) \epsilon_t \pi_{t,k} \right),^1 \quad (1)$$

where we define the event decay function  $f$  to be a simple

<sup>1</sup>Throughout this work, we use a non-traditional parameterization of the gamma distribution. Recall the shape  $a$  and

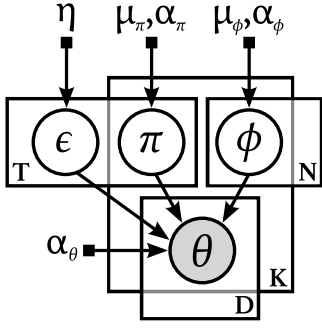


Figure 2: A directed graphical model of Capsule to show considered dependencies. Shaded document topics nodes  $\theta$  are observed. unshaded nodes are hidden variables—these are event occurrences  $\epsilon$ , event content descriptions  $\pi$ , and entity typical concerns  $\phi$ . Plates denote replication: there are  $D$  documents,  $T$  time steps,  $N$  entities, and  $K$  topics. Hyperparameters  $\eta$ ,  $\mu$ , and  $\alpha$  are fixed.

linear decrease:

$$f(t, m) = \begin{cases} 1 - \frac{m-t}{\delta}, & \text{if } t \leq m < t + \delta \\ 0, & \text{otherwise,} \end{cases}$$

with  $\delta$  being the time duration after the event at time  $t$  is no longer relevant. Figure 2 shows the dependencies between the hidden and observed variables as a graphical model.

To complete the specification of all the variables, we place priors on all of the hidden variables. Author concerns  $\phi_{n,k}$  and event content  $\pi_t$  are specified with Gamma priors. Event occurrence  $\epsilon_t$  has a Poisson prior.

**Learning the hidden variables.** In order to use the Capsule model to explore the observed documents, we must compute the posterior distribution. Conditional on the observed document topics  $\theta$ , our goals to compute the posterior values of the hidden parameters—event occurrences  $\epsilon$  and descriptions  $\pi$ , as well as entity concerns  $\phi$ .

As is common for Bayesian models, the exact posterior for Capsule is not tractable to compute; approximating it is our central statistical and computational problem. We develop an approximate inference algorithm for Capsule based on variational methods [38].<sup>2</sup>

Variational inference approaches the problem of posterior inference by minimizing the KL divergence from an approximating distribution  $q$  to the true posterior  $p$ . This is equivalent to maximizing the ELBO:

$$\mathcal{L}(q) = \mathbb{E}_{q(\epsilon, \pi, \phi)} [\log p(\theta, \epsilon, \pi, \phi) - \log q(\epsilon, \pi, \phi)]. \quad (2)$$

scale  $b$  parameterization of the gamma, or

$$\text{Gamma}^*(x | a, b) = \frac{1}{\Gamma(a)b^a} x^{(a-1)} e^{-x/b}.$$

In our alternative parameterization, we use a single mean parameter  $\mu$  and a fixed sparsity hyperparameter  $\alpha$ , or

$$\text{Gamma}(x | \mu, \alpha) = \text{Gamma}^*(x | \alpha, \mu/\alpha);$$

when a gamma distribution is only specified by a single parameter, it is the mean  $\mu$  and the sparsity hyperparameter  $\alpha$  is hidden for simplicity.

<sup>2</sup>Source code available at <https://github.com/ajbc/capsule>. TODO: release on github repo (this link is not active.) submission is not anonymous, so why not?



Figure 3: Screenshots of Capsule visualization of US State Department cables. Left: top words in a topic (manually labeled topic title). Center-top: events over time (height is volume of messages sent, color is probability of an event occurring). Center-bottom: topics for an event on `<date TODO: cyprus coup?>`. Right-top: cyprus entity topics? TODO. Right-bottom: entities shown on a map.

We define the approximating distribution  $q$  using the mean field assumption:

$$q(\epsilon, \pi, \phi) = \prod_{t=1}^T q(\epsilon_t | \lambda_t^\epsilon) \prod_{k=1}^K \left[ \prod_{n=1}^N q(\phi_{n,k} | \lambda_{n,k}^\phi) \prod_{t=1}^T q(\pi_{t,k} | \lambda_{t,k}^\pi) \right]. \quad (3)$$

The variational distributions  $q(\pi)$  and  $q(\phi)$  are both gamma-distributed with free variational parameters  $\lambda^\pi$  and  $\lambda^\phi$ , respectively. The variational distribution  $q(\epsilon)$  is Poisson-distributed with variational parameter  $\lambda^\epsilon$ .

The expectations under  $q$ , which are needed to maximize the ELBO, do not have a simple analytic form, so we use “black box” variational inference techniques [31]. Black box techniques optimize the ELBO directly with stochastic optimization [33]. Full details on our inference algorithm can be found in the appendix. This algorithm produces a fitted variational distribution which can then be used as a proxy for the true posterior, allowing us to explore a collection of documents with Capsule.

**Visualization.** Capsule is a high-level statistical tool. In order to understand and explore its results, a user must scrutinize numerical distributions. To make Capsule more accessible, we developed an open source tool for visualizing its results.<sup>3</sup> Our tool creates a navigator of the documents and latent parameters, allowing users to explore events, entities, topics, and the original documents. Figure 2 shows several screenshots of this browsing interface.

### 3. EVALUATION

In this section we study the performance of Capsule. Using simulated data, we compare Capsule to deterministic methods of event detection and show that Capsule outperforms them at identifying when events occur. We conclude by exploring three real-worlds datasets with Capsule.

#### 3.1 Performance

We generated ten simulated datasets using our generative process. Each dataset spans 100 days and contains content associated with ten entities. Approximately ten events also

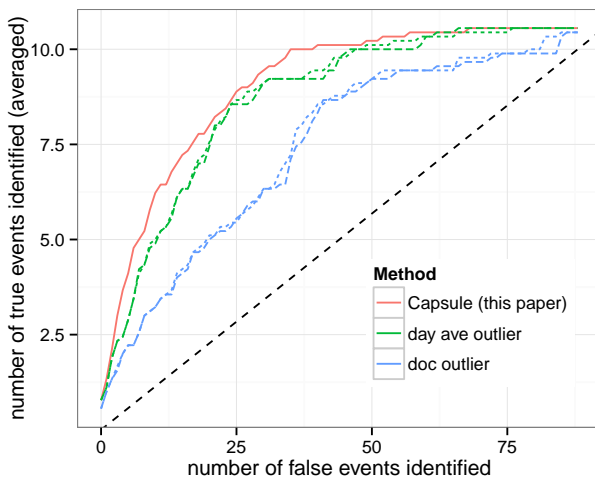
<sup>3</sup>Source code: <https://github.com/ajbc/capsule-viz>. TODO

exist in each dataset, randomly distributed in time and with a three day decay of relevancy.

To evaluate performance, we rank each day by its probability of having an event occur, and plot the number of true events discovered against the number of false positive events, as shown in Figure 4; the area under the curve (AUC) can be computed for a single evaluation metric. Note that this approach is only valid when true events are known, and thus we only apply it to simulated data.

We compare Capsule to two baseline approaches: one considers the greatest document outlier on a given day—days with the furthest outliers are the most likely to have events. The other approach is similar: days are represented by an average of all documents associated with that day, and one considers how these averages deviate from the global average—the further away, the more likely an event.

Figure 4 shows that Capsule outperforms both of these approaches. It should be noted that inference on Capsule will produce different results, depending on the random seed; the results shown are the best of three random seeds.



**Figure 4: Average performance on ten simulated datasets; lines closer to the upper-left are better. Baselines consider outliers based on full corpus averages (dashed) and averages of all entity documents (dotted). Capsule performance is best of three random seeds.**

### 3.2 Exploration

**Cables** ¶ where did we get it / size / preprocessing  
 ¶ plot of events timeline with select real-world match  
 evetns pointed out (verified by history lab)  
 ¶ example interesting entities + figure  
 ¶ explore pairwise entities? (quick with and single figure shared with enron); compare sender vs reiever for same pair (or does direction matter?? tyr both ways) look at sender in norma model vs sender in a few pairs under this construction  
**arXiv** ¶ where did we get it / size / preprocessing  
 ¶ plot of events timeline with select real-world match  
 evetns pointed out (verified by history lab)  
 ¶ example interesting entities + figure  
**enron** ¶ where did we get it / size / preprocessing  
 ¶ plot of events timeline with select real-world match

evetns pointed out (verified by history lab)

- ¶ example interesting entities + figure
- ¶ explore pairwise entities?

## 4. DISCUSSION

We have presented Capsule, a Bayesian model that identifies when events occur, characterizes these events, and discovers the typical concerns of author entities. We have shown that Capsule outperforms deterministic baseline methods and explored its results on three real-world datasets. We anticipate that Capsule and its visualization can be used by historians, political scientist, and others who wish to explore and investigate events in large text corpora. Future work includes expanding the model to incorporate messages recipients and allowing events to impact only a subset of entities.

## 5. REFERENCES

- [1] J. Allan, R. Papka, and V. Lavrenko. On-line new event detection and tracking. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 37–45. ACM, 1998.
- [2] H. Becker, M. Naaman, and L. Gravano. Learning similarity metrics for event identification in social media. In *Proceedings of the third ACM international conference on Web search and data mining*, pages 291–300. ACM, 2010.
- [3] D. M. Blei. Probabilistic topic models. *Communications of the ACM*, 55(4):77–84, 2012.
- [4] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet allocation. *JMLR*, 3:993–1022, Mar. 2003.
- [5] T. Brants, F. Chen, and A. Farahat. A system for new event detection. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*, pages 330–337. ACM, 2003.
- [6] G. Casella and C. P. Robert. Rao-blackwellisation of sampling schemes. *Biometrika*, 83(1):81–94, 1996.
- [7] D. Chakrabarti and K. Punera. Event summarization using tweets. *ICWSM*, 11:66–73, 2011.
- [8] A. J.-B. Chaney and D. M. Blei. Visualizing topic models. In *ICWSM*, 2012.
- [9] K. Das, J. Schneider, and D. B. Neill. Anomaly pattern detection in categorical datasets. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 169–176. ACM, 2008.
- [10] A. Das Sarma, A. Jain, and C. Yu. Dynamic relationship and event discovery. In *Proceedings of the fourth ACM international conference on Web search and data mining*, pages 207–216. ACM, 2011.
- [11] Y. N. Dauphin, H. de Vries, J. Chung, and Y. Bengio. Rmsprop and equilibrated adaptive learning rates for non-convex optimization. *arXiv preprint arXiv:1502.04390*, 2015.
- [12] W. Dou, X. Wang, D. Skau, W. Ribarsky, and M. X. Zhou. Leadline: Interactive visual analysis of text data through event identification and exploration. In *Visual Analytics Science and Technology (VAST), 2012 IEEE Conference on*, pages 93–102. IEEE, 2012.

- [13] G. P. C. Fung, J. X. Yu, P. S. Yu, and H. Lu. Parameter free bursty events detection in text streams. In *Proceedings of the 31st international conference on Very large data bases*, pages 181–192. VLDB Endowment, 2005.
- [14] W. Gao, P. Li, and K. Darwish. Joint topic modeling for event summarization across news and social media streams. In *Proceedings of the 21st ACM international conference on Information and knowledge management*, pages 1173–1182. ACM, 2012.
- [15] Z. Ghahramani and M. J. Beal. Propagation algorithms for variational bayesian learning. *Advances in neural information processing systems*, pages 507–513, 2001.
- [16] V. Guralnik and J. Srivastava. Event detection from time series data. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 33–42. ACM, 1999.
- [17] M. Hoffman, F. R. Bach, and D. M. Blei. Online learning for latent dirichlet allocation. In *advances in neural information processing systems*, pages 856–864, 2010.
- [18] A. Ihler, J. Hutchins, and P. Smyth. Learning to detect events with markov-modulated poisson processes. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(3):13, 2007.
- [19] A. Jackoway, H. Samet, and J. Sankaranarayanan. Identification of live news events using twitter. In *Proceedings of the 3rd ACM SIGSPATIAL International Workshop on Location-Based Social Networks*, pages 25–32. ACM, 2011.
- [20] J. F. C. Kingman. *Poisson Processes*. 1993.
- [21] J. Kleinberg. Bursty and hierarchical structure in streams. *Data Mining and Knowledge Discovery*, 7(4):373–397, 2003.
- [22] G. Kumaran and J. Allan. Text classification and named entities for new event detection. In *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 297–304. ACM, 2004.
- [23] J. H. Lau, N. Collier, and T. Baldwin. On-line trend analysis with topic models: \# twitter trends detection topic model online. In *COLING*, pages 1519–1534, 2012.
- [24] Z. Li, B. Wang, M. Li, and W.-Y. Ma. A probabilistic model for retrospective news event detection. In *Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 106–113. ACM, 2005.
- [25] S. W. Linderman and R. P. Adams. Discovering latent network structure in point process data. *arXiv preprint arXiv:1402.0914*, 2014.
- [26] X. Liu, R. Troncy, and B. Huet. Using social media to identify events. In *Proceedings of the 3rd ACM SIGMM international workshop on Social media*, pages 3–8. ACM, 2011.
- [27] M. Mathioudakis, N. Bansal, and N. Koudas. Identifying, attributing and describing spatial bursts. *Proceedings of the VLDB Endowment*, 3(1-2):1091–1102, 2010.
- [28] A. McCallum, K. Nigam, et al. A comparison of event models for naive bayes text classification. In *AAAI-98 workshop on learning for text categorization*, volume 752, pages 41–48. Citeseer, 1998.
- [29] D. B. Neill, A. W. Moore, M. Sabhnani, and K. Daniel. Detection of emerging space-time clusters. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, pages 218–227. ACM, 2005.
- [30] W. Peng, C. Perng, T. Li, and H. Wang. Event summarization for system management. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1028–1032. ACM, 2007.
- [31] R. Ranganath, S. Gerrish, and D. M. Blei. Black box variational inference. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics, AISTATS '14*, pages 814–822, 2014.
- [32] T. Reuter and P. Cimiano. Event-based classification of social media streams. In *Proceedings of the 2nd ACM International Conference on Multimedia Retrieval*, page 22. ACM, 2012.
- [33] H. Robbins and S. Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [34] T. Sakaki, M. Okazaki, and Y. Matsuo. Earthquake shakes twitter users: real-time event detection by social sensors. In *Proceedings of the 19th international conference on World wide web*, pages 851–860. ACM, 2010.
- [35] H. Sayyadi, M. Hurst, and A. Maykov. Event detection and tracking in social streams. In *ICWSM*, 2009.
- [36] A. Schein, J. Paisley, D. M. Blei, and H. Wallach. Bayesian poisson tensor factorization for inferring multilateral relations from sparse dyadic event counts. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1045–1054. ACM, 2015.
- [37] C. VanDam. A probabilistic topic modeling approach for event detection in social media. Master’s thesis, Michigan State University, 2012.
- [38] M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Found. Trends Mach. Learn.*, 1(1-2):1–305, Jan. 2008.
- [39] X. Wang, C. Zhai, X. Hu, and R. Sproat. Mining correlated bursty topic patterns from coordinated text streams. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 784–793. ACM, 2007.
- [40] G. M. Weiss and H. Hirsh. Learning to predict rare events in event sequences. In *KDD*, pages 359–363, 1998.
- [41] Y. Yang, T. Ault, T. Pierce, and C. W. Lattimer. Improving text categorization methods for event tracking. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval*, pages 65–72. ACM, 2000.
- [42] Y. Zhang, J. Callan, and T. Minka. Novelty and redundancy detection in adaptive filtering. In *Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 81–88. ACM, 2002.
- [43] Q. Zhao, P. Mitra, and B. Chen. Temporal and

information flow based event detection from social text streams. In *AAAI*, volume 7, pages 1501–1506, 2007.

- [44] W. X. Zhao, R. Chen, K. Fan, H. Yan, and X. Li. A novel burst-based text representation model for scalable event detection. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Short Papers-Volume 2*, pages 43–47. Association for Computational Linguistics, 2012.

## APPENDIX

In this appendix, we describe the details of the variational inference algorithm for Capsule. This algorithm fits the parameters of the variational distribution  $q$  in Eq. 3 so that it is close in KL divergence to the posterior.

Recall that the variational distributions  $q(\pi)$  and  $q(\phi)$  are both gamma-distributed with free variational parameters  $\lambda^\pi$  and  $\lambda^\phi$ , respectively. Each parameter  $\lambda$  has two components: sparsity  $\alpha$  and mean  $\mu$ , which parameterize a shape-rate gamma as  $\text{Gamma}(\alpha, \mu/\alpha)$ , as noted previously. Because these parameters are free, we use the softplus function  $\mathcal{P}(x) = \log(1 + \exp(x))$  to constrain them so that they do not violate the requirements of the gamma distribution. The variational distribution  $q(\epsilon)$  is Poisson-distributed with variational parameter  $\lambda^\epsilon$ , which is also constrained by the softplus function.

Minimizing the KL divergence between the true posterior  $p$  and the variational approximation  $q$  is equivalent to maximizing the ELBO (Eq. 2). This maximization is often achieved with closed form coordinate updates, but the Capsule model is not specified with the required conjugate relationships that make this approach possible [15]. Instead, we rely on “black box” variational inference techniques [31] to perform this optimization.

Black box techniques optimize the ELBO directly with stochastic optimization, which maximizes a function using noisy estimates of its gradient [33]. In this case, the function is the ELBO, and we take derivatives with respect to each of the variational parameters. To obtain the noisy estimates, we sample from the variational approximation  $q$ ; these samples then give us the noisy, unbiased gradients used to update our parameters.

It is essential to employ variance reducing techniques; without them, the algorithm would converge too slowly to be of practical value. Details on each of these techniques may be found in the original black box variational inference paper [31].

One of these techniques is Rao-Blackwellization [6]: for each variable, we can write the log probability of all terms containing that variable, giving us

$$\log p_t^\epsilon \triangleq \log p(\epsilon_t | \eta_\epsilon) + \sum_{d \in D_t} \sum_k \log p(\theta_{d,k} | \dots),^4$$

$$\log p_{t,k}^\pi \triangleq \log p(\pi_{t,k} | \mu_\pi, \alpha_\pi) + \mathbf{1}_{\epsilon_t} \sum_{d \in D_t} \log p(\theta_{d,k} | \dots),^5$$

and

$$\log p_{n,k}^\phi \triangleq \log p(\phi_{n,k} | \mu_\phi, \alpha_\phi) + \sum_{d \in D} \log p(\theta_{d,k} | \dots).$$

Then we can write the gradients with respect to the variational parameters as:

$$\nabla_{\lambda_t^\epsilon} \mathcal{L} = \mathbb{E}_q [\nabla_{\lambda_t^\epsilon} \log q_t^\epsilon (\log p_t^\epsilon - \log q_t^\epsilon)],^6$$

$$\nabla_{\lambda_{t,k}^\pi} \mathcal{L} = \mathbb{E}_q [\nabla_{\lambda_{t,k}^\pi} \log q_{t,k}^\pi (\log p_{t,k}^\pi - \log q_{t,k}^\pi)],$$

and

$$\nabla_{\lambda_{n,k}^\phi} \mathcal{L} = \mathbb{E}_q [\nabla_{\lambda_{n,k}^\phi} \log q_{n,k}^\phi (\log p_{n,k}^\phi - \log q_{n,k}^\phi)].$$

Using these gradients, we construct our black box algorithm below in Algorithm 1. As shown, the algorithm does not subsample documents, but for large corpora, we subsample  $B$  documents at each iteration and scale the contribution of these samples by  $D/B$ .

While not shown explicitly in Algorithm 1, we also use control variates and RMSProp [11] to reduce variance.<sup>7</sup> Additionally, we truncate in two instances: sampled gamma variables are given a lower bound to avoid sampling too close to zero, and free parameters are given both lower and upper bounds—the latter is to avoid overflow.

**For Reference** The gamma distribution and derivatives:

$$\begin{aligned} \log \text{Gamma}(x | \mu, \alpha) &= \alpha \log \alpha - \alpha \log \mu - \log \Gamma(\alpha) \\ &\quad + (\alpha - 1) \log x - \frac{\alpha x}{\mu}, \end{aligned} \quad (4)$$

$$\nabla_\mu \log \text{Gamma}(x | \mu, \alpha) = -\frac{\alpha}{\mu} + \frac{\alpha x}{\mu^2}, \quad (5)$$

$$\begin{aligned} \nabla_\alpha \log \text{Gamma}(x | \mu, \alpha) &= \log \alpha + 1 - \log \mu - \Psi(\alpha) \\ &\quad + \log x - \frac{x}{\mu}. \end{aligned} \quad (6)$$

The Poisson distribution and derivative:

$$\log \text{Poisson}(x | \lambda) = x \log \lambda - \log(x!) - \lambda, \quad (7)$$

$$\nabla_\lambda \log \text{Poisson}(x | \lambda) = \frac{x}{\lambda} - 1. \quad (8)$$

<sup>4</sup>Note that we abbreviate

$$p_t^\epsilon = p_t^\epsilon(\theta, \epsilon, \pi, \phi)$$

and

$$p(\theta_{d,k} | \dots) = p(\theta_{d,k} | \epsilon_t, \pi_{t,k}, \phi_{n_{d,k}}, \alpha_\theta),$$

and define

$$D_t \triangleq \forall d \in D : f(t, m_d) \neq 0.$$

<sup>5</sup>We use the indicator shorthand:

$$\mathbf{1}_{\epsilon_t} = \begin{cases} 0, & \text{if } \epsilon_t = 0 \\ 1, & \text{otherwise.} \end{cases}$$

<sup>6</sup>We employ yet another abbreviation:

$$q_t^\epsilon = q(\epsilon_t | \lambda_t^\epsilon).$$

<sup>7</sup>In a conversation with Ranganath, he suggested replacing AdaGrad with RMSprop in setting the learning rate.

The softplus function and derivative:

$$\mathcal{P}(x) = \log(1 + e^x),$$

$$\mathcal{P}'(x) = \frac{e^x}{1 + e^x}. \quad (9)$$

Note that the derivatives in Equations 5, 6, and 8 will always be used in conjunction with Equation 9, as part of the chain rule:

$$\frac{d}{dx} f(\mathcal{P}(x)) = \mathcal{P}'(x) f'(\mathcal{P}(x)). \quad (10)$$

---

**Algorithm 1:** Inference for Cables Model

---

**Input:** document topics  $\theta$   
**Output:** estimates of latent parameters event occurrences  $\epsilon$ , event topics  $\pi$ , and entity topics  $\phi$   
**Initialize**  $\lambda^\epsilon$ ,  $\lambda^\phi$ , and  $\lambda^\pi$  to respective priors  
**Initialize** iteration count  $i = 0$  and  $\sigma^\pi = 0$   
**while** *change in validation likelihood*  $< \Delta$  **do**  
  **for** *each sample*  $s = 1, \dots, S$  **do**  
    **for** *each entity*  $n$  **and** *component*  $k$  **do**  
      draw sample entity topics  
       $\phi_{n,k}[s] \sim \text{Gamma}(\mathcal{P}(\lambda_{n,k}^\phi))$   
      set  $p$ ,  $q$ , and  $g$  using Equations 4–6, 9, and 10:  
       $p_{n,k}^\phi[s] = \log p(\phi_{n,k}[s] | \mu_\phi, \alpha_\phi)$   
       $q_{n,k}^\phi[s] = \log q(\phi_{n,k}[s] | \mathcal{P}(\lambda_{n,k}^\phi))$   
       $g_{n,k}^\phi[s] = \nabla_{\lambda_{n,k}^\phi} \log q(\phi_{n,k}[s] | \mathcal{P}(\lambda_{n,k}^\phi))$   
    **end**  
    **for** *each time step*  $t$  **do**  
      draw sample event occurrence  
       $\epsilon_t[s] \sim \text{Poisson}(\mathcal{P}(\lambda_t^\epsilon))$   
      set  $p$ ,  $q$ , and  $g$  using Equations 7–10:  
       $p_i^\epsilon[s] = \log p(\epsilon_i[s] | \eta)$   
       $q_i^\epsilon[s] = \log q(\epsilon_i[s] | \mathcal{P}(\lambda_i^\epsilon))$   
       $g_i^\pi[s] = \nabla_{\lambda_i^\epsilon} \log q(\epsilon_i[s] | \mathcal{P}(\lambda_i^\epsilon))$   
      **if**  $\epsilon_i[s] \neq 0$  **then**  
        **for** *each component*  $k$  **do**  
          draw sample event topics  
           $\pi_{t,k}[s] \sim \text{Gamma}(\mathcal{P}(\lambda_{t,k}^\pi))$   
          set  $p$ ,  $q$ , and  $g$  using Equations 4–6, 9, and 10:  
           $p_{ik}^\pi[s] = \log p(\pi_{ik}[s] | \alpha_0, \beta_0)$   
           $q_{ik}^\pi[s] = \log q(\pi_{ik}[s] | \lambda_{ik}^\pi)$   
           $g_{ik}^\pi[s] = \nabla_{\lambda_{ik}^\pi} \log q(\pi_{ik}[s] | \lambda_{ik}^\pi)$   
        **end**  
      **end**  
    **end**  
  **end**  
  **for** *each document*  $d$ , *sample*  $s$  **and** *component*  $k$  **do**  
    set  $\mu_{d,k}[s] = \phi_{n_d,k}[s] + \sum_t f(t, m_d) \epsilon_t[s] \pi_{t,k}[s]$   
    set  $p_{d,k}^\theta[s] = \log p(\theta_{d,k} | \mu_{d,k}[s], \alpha_\theta)$  (Eqn. 4)  
     $p_{n,k}^\phi[s] += p_{n,k}^\theta[s]$   
    **for** *each timestep*  $t$  **where**  $t \leq m_d < t + \delta$  **do**  
       $p_t^\epsilon[s] += \sum_k p_{d,k}^\theta[s]$   
      **if**  $\epsilon_t[s] \neq 0$  **then**  
         $p_{t,k}^\pi[s] += p_{t,k}^\theta[s]$   
        update  $\sigma_t^\pi += 1$   
      **end**  
    **end**  
  **end**  
  set  $\hat{\nabla}_{\lambda^\phi} \mathcal{L} \triangleq \frac{1}{S} \sum_s g^\phi[s] (p^\phi[s] - q^\phi[s])$   
  set  $\hat{\nabla}_{\lambda^\epsilon} \mathcal{L} \triangleq \frac{1}{S} \sum_s g^\epsilon[s] (p^\epsilon[s] - q^\epsilon[s])$   
  set  $\hat{\nabla}_{\lambda^\pi} \mathcal{L} \triangleq \frac{1}{\sigma^\pi} \sum_s g^\pi[s] (p^\pi[s] - q^\pi[s])$   
  set  $\rho = (t + \tau)^\kappa$   
  set  $\lambda^\pi += \rho \hat{\nabla}_{\lambda^\pi} \mathcal{L}$   
  set  $\lambda^\epsilon += \rho \hat{\nabla}_{\lambda^\epsilon} \mathcal{L}$   
  set  $\lambda^\phi += \rho \hat{\nabla}_{\lambda^\phi} \mathcal{L}$   
**end**  
  set  $\mathbb{E}[\pi] = \lambda^{\pi,a}$   
  set  $\mathbb{E}[\phi] = \lambda^{\phi,a}$   
  set  $\mathbb{E}[\epsilon] = \lambda^\epsilon$   
**return**  $\mathbb{E}[\pi]$ ,  $\mathbb{E}[\phi]$ ,  $\mathbb{E}[\epsilon]$

---