

Exploring Emerging Issues in Social Torrent via Link-Irregularity Detection

Harish D. Patil, Prof. Ajay kumar Kurra

Abstract— Detection of rising subject matters is now receiving renewed interest encouraged by way of the rapid growth of social networks. Traditional time period frequency based methods might not be right in this context, when you consider that the understanding exchanged in social-community posts comprise no longer simplest textual content but in addition images, URL's and movies. We focus on emergence of topics signaled with the aid of social aspects of these networks. Especially, we focal point on mentions of user links between customers which might be generated dynamically (intentionally or unintentionally) through replies, mentions, and re-tweets. We recommend a probability mannequin of the bringing up conduct of a social society person, and advice to detect the issue of a brand new matter from the anomalies measured via the mannequin. Aggregating anomaly ratings from number of users, we show that we are enable to become aware of emerging themes only founded on the reply/point out relationships in social-network posts. We exhibit our system in several real data sets we gathered from Twitter. The experiments show that the proposed mention anomaly cantered procedures can detect new subject matters at least as early as text-anomaly-situated strategies, and in some instances so much previous when the subject is poorly identified via the textual contents in posts.

Index Terms— Emerging topics, social n/w, anomaly.

I. INTRODUCTION

Communication over social networks sites, such as Face book, Twitter and LinkedIn is gaining its importance in our day by day life. After all the knowledge altering over social networks are not only texts but also URL's images, and videos, they are challenging test beds for the study of information digging. In appropriate, we are excited in the problem of detecting emerging topics from social streams, which can be used to create electronic breaking disclosure, or invent undisclosed market needs or covered economical migration. Compared to traditional media, communal media are capable to capture the original voice of traditional peoples. Therefore, the test is to find the evolution of an issue as soon as possible at a modest no of false positives. The difference that makes social media communal is the continuation of remarks. Here, we mean by remarks links to other users of the same social network in the form of message-to, reply-to, re-tweet-of, or explicitly in the text. One post may contain a number of remarks. Part of end users may include tagged in their posts rarely other users may be tagged their pals all the time. Some users (like celebrities) may receive remarks every minute for others being specified might be an unusual incident. In this sense, remarks are like a

Harish D. Patil, Dept. of Computer Science & Engineering Vathsalya Institute of Science and Technology, Telangana, India.

Prof. Ajaykumar Kurra, Dept. of Computer Science & Engineering Vathsalya Institute of Science and Technology, Telangana, India.

language with the number of words equal to the number of users in a social network. We are excited in finding emerging topics from social network stream based on observing the quoted behavior of peoples. Our basic hypothesis is that a new (emerging) topic is something user feels like conversation, opinion, or sending the information further to their pals. Traditional ways for topic exposure have mainly been concerned with the frequencies of words. A term-density-based approve could suffer from the ambiguity caused by metonyms or homonyms. It may also require knotty pre-processing (for e.g. distribution) depending on the target language. Furthermore, it cannot be exercised when the contents of the messages are mostly non-textual information. On the other hand, the “words” formed by mentions are unique, require little pre processing to gather the information is often separated from the contents, and are available behind hand of the attributes of the contents. Let’s take one simple example and illustrate the transformation of a topic through posts on social networks. The first post by Bob contains mentions to Alice and John, which are both probably friends of Bob, so there is nothing unusual here.



Fig. 1 Overall Example of the emergence of a topic in streams

Fig. 1 Overall Example of the emergence of a topic in social streams The second post by John is a reply to Bob but it is also visible to many friends of John that are not direct friends of Bob. Then in the third post, Dave, one of John’s friends, forwards (called re-tweet in Twitter) the information further down to his own friends. It is worth mentioning that it is not clear what the topic of this discussion is about from the text data, because they are chattering about things (a new cell phones, bikes, cars, or jewellery) that are shown as a link in the text.

In this paper, we recommend a probability model that may capture the typical citing motion of a person, which comprises the quantity of mentions per post and the frequency of users occurring within the mentions. This miniature is used to calculate the oddity of future consumer form. Utilizing the recommended probability model, we will go in to element and measure the individuality or possible blow of a post mirrored within the bringing up behavior of the person. We mixture the ambiguity rankings got on this

manner over 1000's of customers and follow a newly urged change factor detection method headquartered on the always discounting dispensed maximal-possibility coding. This method can catch a change within the statistical dependency design within the time sequence of aggregated anomaly scores, and pinpoint where the topic emergence is see determine 1 The effectiveness of the proposed method is demonstrated on 4 information units now we have collected from Twitter. We exhibit that our point out-anomaly-centred systems can become aware of the emergence of a new subject at least as speedy as text-anomaly-established counterparts.

II. RELATED WORK

Detection and monitoring of subject matters had been studied broadly within the discipline of subject detection and tracking (TDT) [1]. On this context, the foremost venture is to either classify a brand new report into one of the crucial known issues (monitoring) or to notice that it belongs to not one of the recognized classes. Subsequently, temporal constitution of subject matters has been modeled and analyzed by way of dynamic model resolution [4], temporal text mining [5], and factorial hidden Markov units [6]. One other line of research is involved with formalizing the suggestion of "bursts" in a circulate of records. In his seminal paper, Kleinberg modeled bursts utilizing the time varying Poisson procedure with a hidden discrete process that controls the firing price [2]. Just lately, He and Parker developed a physics-encouraged mannequin of bursts centered on the exchange in the momentum of topics [7]. All of the above-stated stories make use of textual content material of the files, however not the social content of the files. The social content (links) has been utilized within the gain knowledge of quotation networks [8]. However, quotation networks are most of the time analyzed in stationary surroundings. The novelty of the present paper lies in focusing on the social content of the records (posts) and in combining this with a change-factor analysis Ease of Use

III. PROPOSED METHOD

The total drift of the proposed approach is proven in Fig. 2. Each step in the float is described within the subsection.

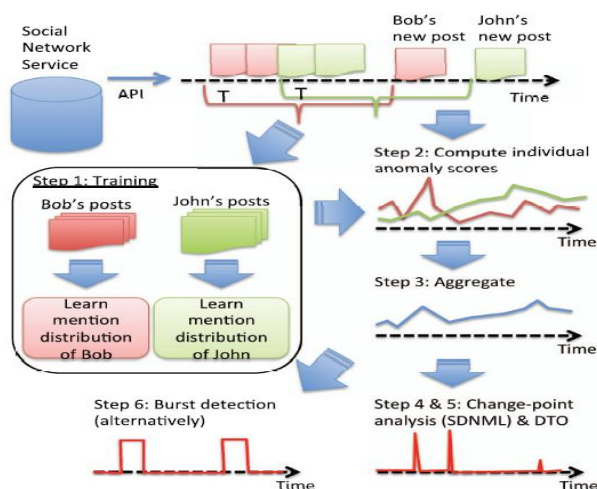


Fig. 2 Overall flow of the proposed method.

We assume that the data arrive from a social community carrier in a sequential manner through some API. For every new submit, we use samples inside the prior time interval of

size T for the corresponding user for training they point out mannequin we advise beneath (Step 1). We assign an anomaly score to each and every submit based on the realized probability distribution (Step 2). The score is then aggregated over customers (Step 3) and extra fed into SDNML-based change point analysis (Steps 4 and 5). We also describe Kleinberg's burst-detection method, which can be used instead of the SDNML-situated change-point evaluation in section 3.6 (Step 6).

3.1 Probability model

On this subsection, we describe the probability model that we used to seize the natural bringing up habits of a person and tips on how to train the model; see Step 1 in Fig. 2. We characterize a publish post in a social network move via the number of mentions k it includes, and the set V of names (IDs) of the mentionees (users who're stated in the publish). There are two types of infinity we have got to take into account here. The primary is the number of customers recounted in a put up. Even though, in apply a user cannot mention 1000's of alternative customers in a submit, we wish to preclude placing an artificial limit on the quantity of users acknowledged in a put up. Rather, we can anticipate a geometrical distribution and combine out the parameter to prevent even an implicit problem via the parameter. The second variety of infinity is the number of users it is easy to almost certainly point out. To restrict limiting the quantity of feasible mentionees, we use a Chinese language Restaurant approach (CRP; see [9]) headquartered estimation; see also Teh et al. [10] who use CRP for limitless vocabulary.

Formally, we keep in mind the next joint likelihood distribution

$$P(k, V | \theta, \{\pi_v\}) = P(k | \theta) \prod_{v \in V} \pi_v. \quad (1)$$

Right here, the joint distribution consists of two ingredients, the probability of the quantity of mentions $k=|V|$ and the chance of each point out given the number of mentions. The chance of the quantity of mentions $P(k|\theta)$ is outlined as a geometric distribution with parameter θ as follows:

$$P(k | \theta) = (1 - \theta)^k \theta. \quad (2)$$

The probability of bringing up user $v \in V$ is denoted by way of π_v (where the sum of π_v over all customers ought to be 1, $\sum_v \pi_v = 1$) and we expect that the k users in V are independently and identically recounted. In different phrases, we ignore the dependences between mentionees and model them as bag of words [11]. Believe that we're given n past posts $T = \{(k1, V1) \dots \dots \dots (kn, Vn)\}g$ from a consumer, and we want to be taught the predictive distribution

$$P(k, V | T) = P(k | T) \prod_{v \in V} P(v | T) \quad (3)$$

The density function of the beta prior distribution is written as follows:

$$p(\theta | \alpha, \beta) = \frac{(1 - \theta)^{\beta-1} \theta^{\alpha-1}}{B(\alpha, \beta)},$$

By way of the Bayes rule, the predictive distribution may also be received as follows:

$$\begin{aligned}
 P(k | T, \alpha, \beta) &= P(k | k_1, \dots, k_n, \alpha, \beta) \\
 &= \frac{P(k, k_1, \dots, k_n | \alpha, \beta)}{P(k_1, \dots, k_n | \alpha, \beta)} \\
 &= \frac{\int_0^1 (1 - \theta) \sum_{i=1}^n k_i + k + \beta - 1 \theta^{n+1+\alpha-1} d\theta}{\int_0^1 (1 - \theta) \sum_{i=1}^n k_i + \beta - 1 \theta^{n+\alpha-1} d\theta}.
 \end{aligned}$$

Each the integrals on the numerator and denominator can be acquired in closed types as beta functions and the predictive distribution may also be rewritten as follows:

$$P(k | T, \alpha, \beta) = \frac{B(n + 1 + \alpha, \sum_{i=1}^n k_i + k + \beta)}{B(n + \alpha, \sum_{i=1}^n k_i + \beta)}.$$

Using the relation between the β function and γ function, we can extra simplify the expression as follows:

$$P(k | T, \alpha, \beta) = \frac{n + \alpha}{m + k + \beta} \prod_{j=0}^k \frac{m + \beta + j}{n + m + \alpha + \beta + j}, \quad (4)$$

Accordingly, the probability of known users is given as follows:

$$P(v | T) = \frac{m_v}{m + \gamma} \quad (\text{for } v: m_v \geq 1). \quad (5)$$

Alternatively, the likelihood of mentioning a brand new person is given as follows:

$$P(\{v : m_v = 0\} | T) = \frac{\gamma}{m + \gamma}. \quad (6)$$

3.2 Computing the link-anomaly score

Step 2 in figure 2, To compute the paradox ranking of a brand new post $x = (t, u, k, V)$ via consumer u at time t containing okay mentions to customers V , we compute the chance (3) with the training set $T_u^{(t)}$, which is the collection of posts via consumer u within the time interval $[t - T, t]$ (we use $T=30$ days in this paper). As a consequence, the link-anomaly ranking is defined as follows:

$$\begin{aligned}
 s(x) &= -\log \left(P(k | T_u^{(t)}) \prod_{v \in V} P(v | T_u^{(t)}) \right) \\
 &= -\log P(k | T_u^{(t)}) - \sum \log P(v | T_u^{(t)}).
 \end{aligned} \quad (7)$$

3.3 Combining anomaly scores from different users

On this subsection, we describe mix the paradox ratings from exclusive users, see Step 3 in figure 2. Paradox rating in (7) is computed for each person relying on the current publish of user u and his/her prior conduct $T_u^{(t)}$. To measure the general trend of consumer conduct, we recommend to mixture the paradox ratings got for posts (x_1, \dots, x_n) making use of a discretization of window size $\tau > 0$ as follows:

$$s_j' = \frac{1}{\tau} \sum_{t_i \in [\tau(j-1), \tau j]} s(x_i), \quad (8)$$

3.4 Change-point detection via sdmml coding

Algorithmically, the exchange-point detection system can also be outlined as follows: For comfort, we denote the aggregated anomaly rating as x_j as a substitute of s_j' .

1. *First-layer learning.* Let $x^{j-1} := \{x_1, \dots, x_{j-1}\}$ be the collection of aggregated anomaly scores from discrete time 1 to $j - 1$. Sequentially learn the SDNML density function $p_{\text{SDNML}}(x_j | x^{j-1}) (j = 1, 2, \dots)$; see Section 3.7 for details.
2. *First-layer scoring.* Compute the intermediate change-point score by smoothing the log loss of the SDNML density function with window size κ as follows:

$$y_j = \frac{1}{\kappa} \sum_{j=j-\kappa+1}^j (-\log p_{\text{SDNML}}(x_j | x^{j-1})).$$

3. *Second-layer learning.* Let $y^{j-1} := \{y_1, \dots, y_{j-1}\}$ be the collection of smoothed change-point score obtained as above. Sequentially learn the second layer SDNML density function $p_{\text{SDNML}}(y_j | y^{j-1}) (j = 1, 2, \dots)$; see Section 3.7 for details.
4. *Second-layer scoring.* Compute the final change-point score by smoothing the log loss of the SDNML density function as follows:

$$\text{Score}(y_j) = \frac{1}{\kappa} \sum_{j=j-\kappa+1}^j (-\log p_{\text{SDNML}}(y_j | y^{j-1})). \quad (9)$$

3.5 Dynamic threshold optimization (dto)

Algorithm 1. Dynamic Threshold Optimization (DTO) [19].

Given: $\{Score_j | j = 1, 2, \dots\}$: scores, N_H : total number of cells, ρ : parameter for threshold, λ_H : estimation parameter, r_H : discounting parameter, M : data size
Initialization: Let $q_1^{(1)}(h)$ (a weighted sufficient statistics) be a uniform distribution.

for $j = 1, \dots, M - 1$ **do**

Threshold optimization: Let l be the least index such that $\sum_{h=1}^l q^{(j)}(h) \geq 1 - \rho$. The threshold at time j is given as

$$\eta(j) = a + \frac{b - a}{N_H - 2} (l + 1).$$

Alarm output: Raise an alarm if $Score_j \geq \eta(j)$.

Histogram update:

$$q_1^{(j+1)}(h) = \begin{cases} (1 - r_H)q_1^{(j)}(h) + r_H & \text{if } Score_j \text{ falls} \\ & \text{into the } h\text{th} \\ & \text{cell,} \\ (1 - r_H)q_1^{(j)}(h) & \text{otherwise.} \end{cases}$$

$$q^{(j+1)}(h) = (q_1^{(j+1)}(h) + \lambda_H) / (\sum_h q_1^{(j+1)}(h) + N_H \lambda_H).$$

end for

3.6 Kleinberg's burst-detection method

Furthermore to the change-point detection situated on SDNML adopted by using DTO described in prior sections, we also scan the mixture of our process with Kleinberg’s burst-detection approach [2]. More particularly, we implemented a two-state variant of Kleinberg’s burst detection model. The purpose we selected the 2-state variant used to be seeing that on this scan we anticipate no hierarchical constitution. The burst-detection method is established on a probabilistic automaton model with two states, burst state and non-burst state. Some routine (e.g., arrival of posts) are assumed to happen in step with a time-various Poisson strategies whose fee parameter depends upon the current state. The burst-detection procedure estimates the state transition sequence $i_t \in \{nonburst, burst\}$ ($t = 1, \dots, n$) that maximizes the likelihood

$$p_{sw}^b (1 - p_{sw})^{n-b} \prod_{t=1} f_{exp}(x_t; \alpha_{i_t}),$$

Where P_{sw} is a given state transition likelihood, b is the quantity of state transitions within the sequence $i_t = (t = 1, \dots, n)$, $f_{exp}(x; \alpha)$ is the chance density function of the exponential distribution with fee parameter α and x_t is the t^{th} interevent interval. The most desirable sequence can be efficiently obtained by means of dynamic programming [2]. To obtain the event occasions and their intervals, we define an occasion as a point in time when the aggregated link anomaly score (8) exceeds a threshold θ_{burst} .

IV. PROPOSED METHOD

4.1 Observation setup

We collected data sets from Twitter. This set is associated with a list of posts in a service called together. Together is a collaborative service where people can tag Twitter posts that are related to each other and organize a List of posts that belong to a certain topic. Our goal is to evaluate whether the proposed approach can detect the Emergence of the topics recognized and collected by people.

TABLE 1

Parameter Values We Used in the Real Data Experiments

Model	Parameter name	value
Mention model	Beta distribution	$\alpha = \beta = 0.5$
	CRP parameter	$\gamma = 0.5$
	Training period	$T = 30$ days
SDNML-based change-point detection	AR model order	$p = 30$
	Smoothing parameter	$\kappa = 15$
Dynamic threshold optimization (DTO)	Number of bins	$N_H = 20$
	Smoothing parameter	$\lambda_H = 0.01$
	Discount rate	$r = 0.005$
	Significance level parameter	$\rho = 0.05$
	Normal data upper limit	$a = \text{average} + 3\sigma$ of the input data
Normal data lower limit	$b = \text{minimum}$ of the input data	
Kleinberg’s burst detection model	Rate parameter (nonburst state)	$\alpha_{nonburst} = 0.001$ (1/s)
	Rate parameter (burst state)	$\alpha_{burst} = 0.01$ (1/s)
	State transition probability	$p_{sw} = 0.3$
	Threshold parameter	$\theta_{burst} = 0.9995$ -quantile point of the aggregated anomaly scores

TABLE 2
Data Sets

Data Set	No Of Participants	No Of Posts	Keywords	Keywords(Unicode)
YouTube	160	5,47,287	Senkaku	u'\u5c16\u95a3'
"NASA"	90	2,78,115	Arsenic	u'\u30d2\u7820'

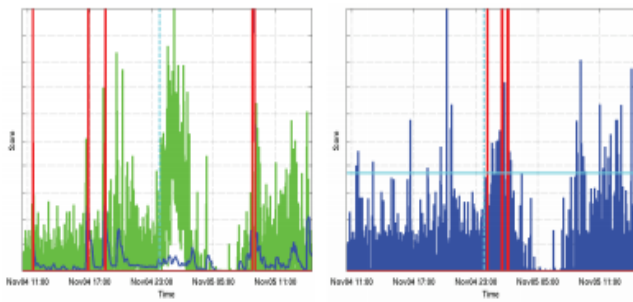
The information units we accumulated is called “NASA” and it’s corresponds to a user prepared record in Together. For all record, we extracted a list of Twitter users that regarded in the list, and picked up Twitter posts from those customers. See table 2 for the quantity of participants and the quantity of posts we amassed for every information set. Word that we amassed Twitter posts up to 30 days earlier than the time interval of interest for every consumer; thus, the number of posts we analyzed used to be much larger than the number of posts listed in Together.

4.2 “YouTube” data set

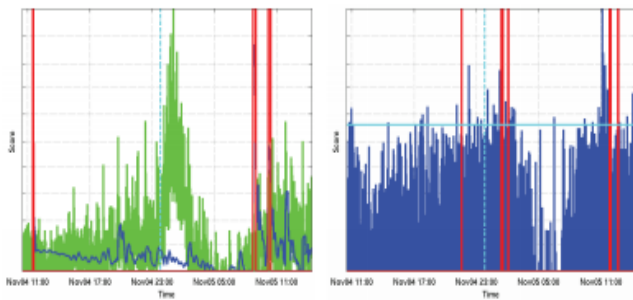
This knowledge set is concerning the contemporary leakage of some exclusive video through the Japan Coastal protect officer. The keyword used in the key phrase-based ways used to be “Senkaku.” Figs. 3a and 3b show the results of link-anomaly based change detection and burst detection, respectively. Figs. 3c and 3d exhibit the results of textual content-anomaly-founded exchange detection and burst detection, respectively. Figs. 3e and 3f exhibit the results of keyword-frequency-founded exchange detection and burst detection, respectively. The first alarm instances of the proposed hyperlink-anomaly-based change-point analysis and the text-anomaly-centred change-point analysis have been each 08:44, Nov. 05, which have been just about 9 hours after the first put up in regards to the video leakage. Even though there is an elevation within the aggregated anomaly rating (8) in Fig. 3a around midnight, Nov 05, it seems that SDNML fails to become aware of this elevation as a transformation-factor. In fact, the link-anomaly-established burst detection (Fig. 3b) raised an alarm at 00:07, which is prior than the keyword-frequency-established dynamic thresholding and toward the key phrase-frequency based burst detection at 23:59, Nov 04. The alarm time of the textual content-anomaly-based burst detection was once 01:24, Nov 05.

4.3 “NASA” data set

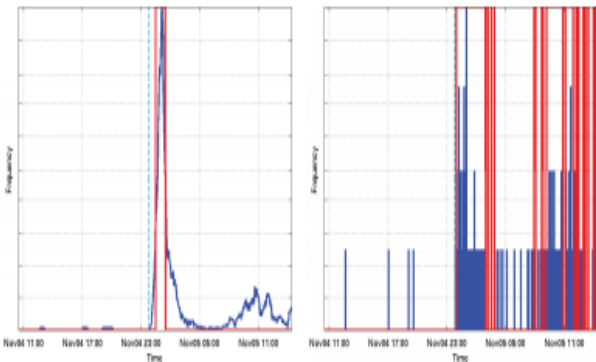
This knowledge set is regarding the dialogue amongst Twitter users serious about astronomy that the key phrase used in the key phrase-based models used to be “arsenic.” Figs. 4a and 4b show the results of link anomaly-centered trade detection and burst detection, respectively. Figs. 4c and 4d exhibit the outcome of text anomaly-established exchange detection and burst detection, respectively. Figs. 4e and 4f show the identical outcome for the key phrase-frequency based ways. The primary alarm times of the two link-anomaly-established approaches have been 22:20,



(a) Link-anomaly-based change-point (b) Link-anomaly-based burst de-analysis. Green: Aggregated anomaly detection. Blue: Aggregated anomaly score (8) at $\tau = 1$ minute. Blue: score (8) at $\tau = 1$ second. Horizontal cyan: threshold for the filtering step in Kleinberg's burst model. Red: Alarm time. Red: Burst state.



(c) Text-anomaly-based change-point (d) Text-anomaly-based burst de-analysis. Green: Aggregated anomaly detection. Blue: Aggregated anomaly score (8) at $\tau = 1$ minute. Blue: score (8) at $\tau = 1$ second. Horizontal cyan: threshold for the filtering step in Kleinberg's burst model. Red: Alarm time. Red: Burst state.

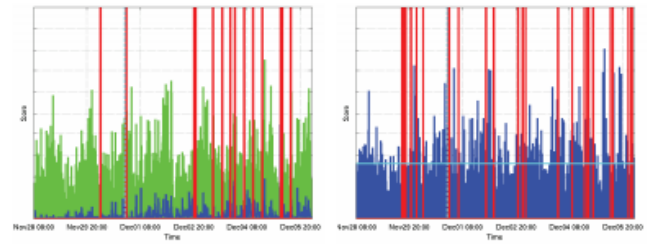


(e) Keyword-frequency-based change-point analysis. Blue: detection. Blue: Frequency of keyword "Senkaku" per one second. Red: Alarm time. Red: Burst state (burst or not).

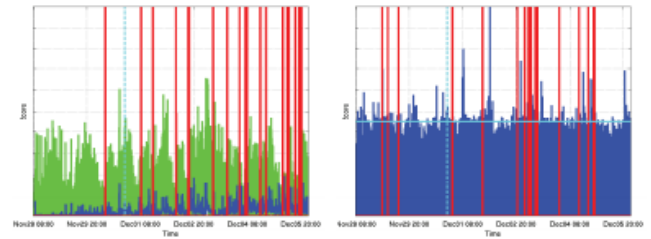
Fig. 3. Result of "YouTube" data set. The first post about the video leakage was posted at 23:48, Nov 04 (indicated by vertical cyan dashed lines).

Nov 30 (alternate-point detection) and 22:44, Nov 30 (burst detection), respectively.

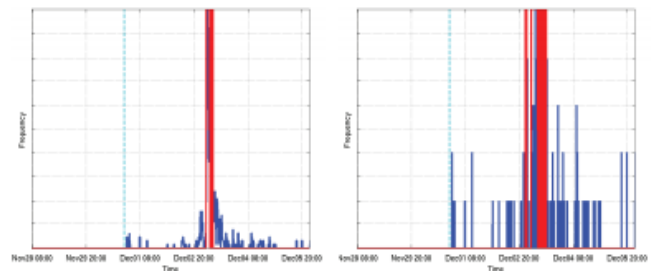
Both of those had been previous than NASA's legit press conference (04:00, Dec 03) and were earlier than the text-anomaly-founded approaches (alternate-point detection at 08:17, Dec 01 and burst detection at 00:54, Dec 01). The keyword-founded methods are even slower.



(a) Link-anomaly-based change-point (b) Link-anomaly-based burst de-analysis. Green: Aggregated anomaly detection. Blue: Aggregated anomaly score (8) at $\tau = 1$ minute. Blue: score (8) at $\tau = 1$ second. Horizontal cyan: threshold for the filtering step in Kleinberg's burst model. Red: Alarm time. Red: Burst state.



(c) Text-anomaly-based change-point (d) Text-anomaly-based burst de-analysis. Green: Aggregated anomaly detection. Blue: Aggregated anomaly score (8) at $\tau = 1$ minute. Blue: score (8) at $\tau = 1$ second. Horizontal cyan: threshold for the filtering step in Kleinberg's burst model. Red: Alarm time. Red: Burst state.



(e) Keyword-frequency-based change-point analysis. Blue: detection. Blue: Frequency of keyword "arsenic" per one second. Red: Alarm time. Red: Burst state (burst or not).

Fig. 4. Result of "NASA" data set. The initial post predicting NASA's finding about arsenic-eating organism was posted at 21:39, Nov 30 much earlier than NASA's official press conference at 04:00, Dec 03.

TABLE 3
Detection Time and the Number of False Detections

Method	"Job hunting"	"Youtube"	"NASA"
		22:50, Jan 08	23:48, Nov 04
Link-anomaly-based change-point detection	3	3	13
	22:55, Jan 08	08:44, Nov 05	22:20, Nov 30
Link-anomaly-based burst detection	2	2	26
	23:04, Jan 08	00:07, Nov 05	22:44, Nov 30
Text-anomaly-based change-point detection	4	2	20
	22:51, Jan 08	08:44, Nov 05	08:17, Dec 01
Text-anomaly-based burst detection	0	4	18
	22:51, Jan 08	01:24, Nov 05	00:54, Dec 01
Keyword-frequency-based dynamic thresholding	0	0	5
	22:57, Jan 08	00:30, Nov 05	04:10, Dec 03
Keyword-frequency-based burst detection	5	14	10
	22:50, Jan 08	23:59, Nov 04	23:59, Dec 02

The key phrase-frequency-centered dynamic thresholding raised an alarm at 04:10, Dec 03 after NASA's official press unencumbered burst detection raised an alarm at 23:59, Dec 02; see Table 3.

V. CONCLUSION

On this paper, we have proposed a new procedure to become aware of the emergence of topics in a social community circulation. The fundamental inspiration of our procedure is to centre of attention on the social aspect of the posts reflected within the bringing up behavior of customers instead of the textual contents. We've proposed a likelihood model that captures both the quantity of mentions per put up and the frequency of mentionee. We've got combined the proposed mention model with the SDNML alternate-factor detection algorithm [3] and Kleinberg's burst-detection model [2] to pinpoint the emergence of a subject. Because the proposed procedure does no longer rely on the textual contents of social community posts, it's effective to rephrasing and it may be applied to the case where issues are involved with knowledge rather than texts, akin to pictures, video, audio, and many others.

We've got utilized the proposed method to 2 real information sets that we've got collected from Twitter. The 2 data sets included a wide-unfold dialogue a few fast propagation of news about a video leaked on YouTube ("YouTube" knowledge set), a rumor about the upcoming press conference via NASA ("NASA" information set). In all the data sets, our proposed strategy showed promising performance. In one out of two information units, the detection by the proposed link-anomaly founded ways used to be earlier than the textual content-anomaly-based counterparts. In addition, for "NASA" data set, wherein the key phrase that defines the subject is extra ambiguous than the primary data set, the proposed hyperlink-anomaly-founded methods have detected the emergence of the topics even prior than the keyword-founded systems that use hand-chosen key phrases. All of the analysis offered on this paper was once performed offline, but the framework itself may also be utilized online. We are planning to scale up the proposed procedure to control social streams in real time. It will even be intriguing to combine the proposed link-anomaly model with text-founded strategies, considering the fact that the proposed hyperlink-anomaly mannequin does now not right away tell what the ambiguity is. Combo of the phrase-established method with the link-anomaly model would advantage each from the efficiency of the point out model and the intuitiveness of the phrase-situated approach.

ACKNOWLEDGMENTS

I feel great in expressing our deepest sense of gratitude to our guide and HOD Prof. Ajaykumar Kurra for his encouragement and enlightened comments throughout this project work. His appreciative suggestion always motivated us for putting most willing efforts on my study during project report. We are also thankful to the concerned authorities who directly or indirectly helped us in the project.

REFERENCES

- [1] J. Allan et al., "Topic Detection and Tracking Pilot Study: Final Report," Proc. DARPA Broadcast News Transcription and Understanding Workshop, 1998.
- [2] J. Kleinberg, "Bursty and Hierarchical Structure in Streams," Data Mining Knowledge Discovery, vol. 7, no. 4, pp. 373-397, 2003.
- [3] Y. Urabe, K. Yamanishi, R. Tomioka, and H. Iwai, "Real-Time Change-Point Detection Using Sequentially Discounting Normalized Maximum Likelihood Coding," Proc. 15th Pacific-Asia Conf. Advances in Knowledge Discovery and Data Mining (PAKDD' 11), 2011.

- [4] S. Morinaga and K. Yamanishi, "Tracking Dynamics of Topic Trends Using a Finite Mixture Model," Proc. 10th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining, pp. 811-816, 2004.
- [5] Q. Mei and C. Zhai, "Discovering Evolutionary Theme Patterns from Text: An Exploration of Temporal Text Mining," Proc. 11th ACM SIGKDD Int'l Conf. Knowledge Discovery in Data Mining, pp. 198-207, 2005.
- [6] A. Krause, J. Leskovec, and C. Guestrin, "Data Association for Topic Intensity Tracking," Proc. 23rd Int'l Conf. Machine Learning (ICML' 06), pp.497-504,2006.
- [7] D. He and D.S. Parker, "Topic Dynamics: An Alternative Model of Bursts in Streams of Topics," Proc. 16th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining, pp. 443-452, 2010.
- [8] H. Small, "Visualizing Science by Citation Mapping," J. Am. Soc. Information Science, vol. 50, no. 9, pp. 799-813, 1999.
- [9] D. Aldous, "Exchangeability and Related Topics," Ecole d' Ete ' de Probabilite ' s de Saint-Flour XIII—1983, pp. 1-198, Springer, 1985.[10] Y. Teh, M. Jordan, M. Beal, and D. Blei, "Hierarchical Dirichlet Processes," J. Am. Statistical Assoc., vol. 101, no. 476, pp. 1566-1581, 2006.
- [11] D. Lewis, "Naive (Bayes) at Forty: The Independence Assumption in Information Retrieval," Proc. 10th European Conf. Machine Learning (ECML'98),pp.4-15,1998.
- [12] K. Yamanishi and J. Takeuchi, "A Unifying Framework for Detecting Outliers and Change Points from non-Stationary Time Series Data," Proc. Eighth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining, 2002.
- [13] J. Takeuchi and K. Yamanishi, "A Unifying Framework for Detecting Outliers and Change Points from Time Series," IEEE Trans. Knowledge Data Eng., vol. 18, no. 4, pp. 482-492, Apr.2006.
- [14] J. Rissanen, "Strong Optimality of the Normalized ML Models as Universal Codes and Information in Data," IEEE Trans. Information Theory, vol.47, no.5, pp.1712-1717, July2001.
- [15] T. Roos and J. Rissanen, "On Sequentially Normalized Maximum Likelihood Models," Proc. Workshop Information Theoretic Methods in Science and Eng., 2008
- [16] J. Rissanen, T. Roos, and P. Myllymäki, "Model Selection by Sequentially Normalized Least Squares," J. Multivariate Analysis, vol. 101, no. 4, pp. 839-849, 2010.
- [17] C. Giurcaneanu, S. Razavi, and A. Liski, "Variable Selection in Linear Regression: Several Approaches Based on Normalized Maximum Likelihood," Signal Processing, vol. 91, pp. 1671-1692, 2011.
- [18] C. Giurcaneanu and S. Razavi, "AR Order Selection in the Case When the Model Parameters Are Estimated by Forgetting Factor Least-Squares Algorithms," Signal Processing, vol. 90, no. 2, pp. 451-466, 2010.
- [19] K. Yamanishi and Y. Maruyama, "Dynamic Syslog Mining for Network Failure Monitoring," Proc. 11th ACM SIGKDD Int'l Conf. Knowledge Discovery in Data Mining, pp. 499-508, 2005.



Harish D. Patil received the BE Degree in Computer from the University Of Pune in 2013 He is working toward the master's degree in the Department of Computer Science and Engineering, from JNTU Hyderabad, India. Interest in Social media and Cloud Computing



Mr. Ajaykumar Kurra has completed B.TECH (CSIT) from JNTU Hyderabad, and M.TECH (CSE) From JNTU Hyderabad. He has 6 years of experience in Academic, Currently working as HOD OF CSE at VATHSALYA INSTITUTE OF SCIENCE AND TECHNOLOGY. Research areas include Data Mining and Wireless Mobile Ad-hoc Networks and Natural Language Processing. He Published 6 international journals.