

Forecasting Protests by Detecting Future Time Mentions in News and Social Media

ABSTRACT

Civil unrest (protests, strikes, and "occupy" events) is a common occurrence in both democracies and authoritarian regimes. The study of civil unrest is a key topic for political scientists as it helps capture an important mechanism by which citizenry express themselves. In countries where civil unrest is lawful, qualitative analysis has revealed that more than 75% of the protests are planned, organized, and/or announced in advance; therefore detecting future time mentions in relevant news and social media is a simple way to develop a protest forecasting system. We develop such a system in this paper, using a combination of key phrase learning to identify what to look for, probabilistic soft logic to reason about location occurrences in extracted results, and time normalization to resolve future tense mentions. We illustrate the application of our system to 10 countries in Latin America, viz. Argentina, Brazil, Chile, Colombia, Ecuador, El Salvador, Mexico, Paraguay, Uruguay, and Venezuela. Results demonstrate our successes in capturing significant societal unrest in these countries with an average lead time of 4.08 days. We also study the selective superiorities of news media versus social media (Twitter, Facebook) and identify relevant tradeoffs.

Categories and Subject Descriptors

H.4 [Information Systems Applications]: Miscellaneous;
D.2.8 [Software Engineering]: Metrics—*complexity measures, performance measures*

General Terms

Theory

Keywords

ACM proceedings, L^AT_EX, text tagging

1. INTRODUCTION

Permission to make digital or hard copies of all or part 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. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

[Sathappan writes: *Dr Mares's text*] Protest is a means by which citizens communicate their views and preferences to those in authority. Representative government is based upon communication between governed and governors; ideally, the channels for communication are open, transparent, credible and efficient. Governments, nevertheless, find it difficult to know on any one issue and at any one time how their constituencies value the available options. Elections are retrospective indicators and rarely issue specific; polling taps into sentiment, but is not a good indicator of priorities or strength of feeling because of the low cost associated with responding. Events, on the other hand, indicate a willingness to bear some costs (organization, mobilization, identification) in support of an issue and thus reveal not only preferences but provide some indication of priorities.

Protest is especially important in democracies that are struggling to consolidate themselves, such as those in Latin America. The combination of weak channels of communication between citizen and government, and a citizenry that still has not grasped the desirability of elections as the means to affect politics means that public protest will be an especially attractive option. To illustrate the power of protest in Latin America we need only recall that between 1985 and 2011 17 Presidents resigned or were impeached under pressure from demonstrations, usually violent, in the streets. Protests have also resulted in the rollback of prices increases for public services, such as in Brazil in June 2013.

We can hypothesize the protests that are larger will be more disruptive and communicate support for its cause better than smaller protests. Mobilizing large numbers of people is more likely to occur if a protest is organized and the time and place publicized. Because protest is costly and more likely to succeed if it is large we should expect planned, rather than spontaneous, protests to be the norm. Indeed, in a sample of 288 events from our study selected for qualitative review of their antecedents, for 225 we located communications regarding the upcoming occurrence of the event and only 49 were classified as 'spontaneous' (we could not determine whether communications had or had not occurred in 14 cases).

2. RELATED WORKS

Three categories of related work – *Event Detection, Extraction of Planned Events and Event Forecasting* – are briefly discussed here.

First, Event Detection/Extraction from textual News has been studied extensively in literature. [3] [57][24] make use of document clustering techniques to identify events retro-

spectively or as the stories arrive.[20],[8], [42] talk about extraction patterns/templates to extract information from text. [45] shows it's possible to accurately extract a calendar of significant events from Twitter by training a tagger for recognizing event phrases.[47] captures tweet clusters of interest to identify late breaking News from twitter.In an altogether different application [46] observes tweets to enable detection of occurrences of Earth Quakes promptly.

Second, some work has been done regarding extraction of planned/future mentions of events from Social Media. RecordedFuture[52] [Sathappan writes: *TODO: write how it differs from our application*] is an analytics company that performs real-time analysis of news and tweets to identify mentions of future events.[51] and [54] use classification and regression techniques to identify the time to an event referred to by a tweet.In [51] a tweet was classified into one of the several identified equal length time bins. [29] tries to provide a collective image of the future associated with an entity summarizing all future related information available.In [10] and [9] content about known planned events is identified from Social Media. Several work has also been done on temporal information extraction from text. [32] presents a search engine ChronoSeeker for searching future and past events.It makes use of an SVM Classifier to disambiguate between the various temporal expressions in a document. [7] and [22] also try to extract future temporal references from text, with the latter using a classifier approach to differentiate between a planned event and a rumor.

Third, few work has been done in the area of Event Forecasting. In [41] the authors learn event sequences from a corpora spanning over 22 years and then use these sequences to say if an event of interest (disease outbreaks, deaths and riots etc) will occur sometime in the future.. In [30], the author makes use of data from RecordedFuture[52] to find if a significant protest event will occur in the subsequent three days using a random forest classifier.The author only focuses on prediction of significant events and also the forecast is limited to the next three days.[21] and [56] are two pieces of research that are very close to our line of work. Both papers follow similar methodologies but are based on different datasets—Twitter and Tumblr.In [21], a list of 335 keywords identified by experts is used to filter twitter stream and the filtered twitter streams are then searched for the presence of future dates in a naive manner by first searching for month names and then for a number less than 31. Such an approach, will not be capable of finding relative mentions of future dates like "tomorrow", "next tuesday" etc. Any location mentioned in the tweet text is used as the event location. If there are no location mentions then the location is determined based on [53].[56] works on Tumblr and makes uses of a fewer set of keywords(59) to filter the Tumblr feed. The filtered feed is further refined by searching for mentions of around 1022 different location names.Finally, the documents are searched for a future date in the same way as in [21].

3. PROBLEM

The problem we try to solve here is the identification of 'Calls for Protest/Strike or any Civil Disobedience movements' in Social Media like News, Blogs, Twitter, Facebook, etc. and to predict the date of event and event location upto a city level resolution.An accurately identified 'call for

protest' is then sent out to an alerting system. An alert is structured as shown in Fig. 8. It is a structured record containing When/Where/Why/Who of the protest and the current date or date at which a forecast is made. The 'when' is specified in granularities of days. The **where** provides a tiered description specifying the (country, state, city), e.g., (Honduras, Francisco Morazan, Tegucigalpa). The **why** (or event type) captures the main objective or reason for a civil unrest event, and is meant to come from 7 broad classes (e.g., 'Employment & Wages', 'Housing', 'Energy & Resources' etc.) each of which is further categorized into whether the event is forecast to be violent or not. Finally, the **who** (or population) denotes common categories of human populations used in event coding [48] such as Business, Ethnic, Legal (e.g. judges or lawyers), Education (e.g. teachers or students or parents of students), Religious (e.g. clergy), Medical (e.g., doctors or nurses), Media, Labor, Refugees/Displaced, Agricultural (e.g. farmers, or just General Population.

Concomitant with the definitions in the above section, a GSR event contains again the where/why/when/who of a protest that has actually occurred and a *reported date* (the date a newspaper reports the protest as having happened). See Fig. 8 (right). The GSR is organized by an independent third party (MITRE) and the authors of this study do not have any participation in this activity.

4. PROBABILISTIC SOFT LOGIC

In this section, we briefly describe Probabilistic Soft Logic(PSL) [33] which is used to geo-code news/blogs as given in 5.1.1.

PSL is a framework for collective probabilistic reasoning on relational domains.PSL models have been developed in various domains, including collective classification [18], ontology alignment [17], personalized medicine [6], opinion diffusion [5], trust in social networks [28], and graph summarization [38].PSL represents the domain of interest as logical atoms. It uses first order logic rules to capture the dependency structure of the domain, based on which it builds a joint probabilistic model over all atoms.Instead of hard truth values of 0 (false) and 1 (true), PSL uses soft truth values relaxing the truth vlaues to the interval [0, 1].The logical connectives are adapted accordingly. This makes it easy to incorporate similarity or distance functions.

User defined *predicates* are used to encode the relationships and attributes and *rules* capture the dependencies and constraints. Each rule's antecedent is a conjunction of atoms and its consequent is a dis-junction. The rules can also labeled with non negative weights which are used during the inference process. The set of predicates and weighted rules thus make up a PSL program where known truth values of ground atoms derived from observed data and unknown truth values for the remaining atoms are learned using the PSL inference.

Given a set of atoms $\ell = \{\ell_1, \dots, \ell_n\}$, an interpretation defined as $I : \ell \rightarrow [0, 1]^n$ is a mapping from atoms to soft truth values. PSL defines a probability distribution over all such interpretaions such that those that satisfy more ground rules are more probable. *Lukasiewicz t-norm* and its corresponding co-norm are used for defining relaxations of the logical AND and OR respectively to determine the degree to which a ground rule is satisfied. Given an interpretation I , PSL defines the formulas for the relaxation of the logical conjunction (\wedge), disjunction (\vee), and negation (\neg) as

Table 1: comparison of our approach with other future event detection methods

	Domain Specific	Multi-Source	Geo-Coding	Temporal Normalization	Feature 4
[32, 54]				✓	
[56]	✓		✓		
reference set 3		✓			
our method	✓	✓	✓	✓	✓

follows:

$$\begin{aligned}\ell_1 \tilde{\wedge} \ell_2 &= \max\{0, I(\ell_1) + I(\ell_2) - 1\}, \\ \ell_1 \tilde{\vee} \ell_2 &= \min\{I(\ell_1) + I(\ell_2), 1\}, \\ \neg \ell_1 &= 1 - I(\ell_1),\end{aligned}$$

The interpretation I determines whether the rules is satisfied, if not, the *distance to satisfaction*. A rule $r \equiv r_{body} \rightarrow r_{head}$ is satisfied if and only if the truth value of head is atleast that of the body. The rule's distance to satisfaction measures the degree to which this condition is violated.

$$d_r(I) = \max\{0, I(r_{body}) - I(r_{head})\}$$

PSL then induces a probability distribution over possible interpretations I over the given set of ground atoms l in the domain. If R is the set of all ground rules that are instances of a rule from the system and uses only the atoms in I then, the probability density function f over I is defined as

$$f(I) = \frac{1}{Z} \exp\left[-\sum_{r \in R} \lambda_r (d_r(I))^p\right] \quad (1)$$

$$Z = \int_I \exp\left[-\sum_{r \in R} \lambda_r (d_r(I))^p\right] \quad (2)$$

where λ_r is the weight of the rule r , Z is the continuous version of the normalization constant used in discrete Markov random fields, and $p \in \{1, 2\}$ provides a choice between two different loss functions, linear and quadratic. The values of the atoms can be further restricted by providing linear equality and inequality constraints allowing one to encode functional constraints from the domain. PSL provides for two kinds of inferences (a)most probable explanation and (b)calculation of the marginal distributions. In the MPE inference given a partial interpretation with grounded atoms based on observed evidence, the PSL program infers the truth values for the unobserved atoms satisfying the most likely interpretation. In the second setting, given ground truth data for all atoms we can learn the weights for the rules in our PSL program.

5. LINGUISTIC PREPROCESSING

As part of the general streaming architecture of the EMBERS system, all textual input (e.g., tweets, news articles, blog postings) is subjected to shallow linguistic processing prior to analysis. Our data set is multilingual, with Spanish, Portuguese and English predominating in our data set. Commercial tools¹ are used for language identification, tokenization, lematization and named entity extraction. The lematized

¹BASIS Technology's Rossette Linguistic Platform []

Table 2: EMBERS system statistics

Archived data	12.4 TB
Archive size	ca. 3 billion messages
Data throughput	200-2000 messages/sec
Daily ingest	15 GB
System memory	50 GB
System core	16 vCPUs
System output	ca. 40 warnings/day

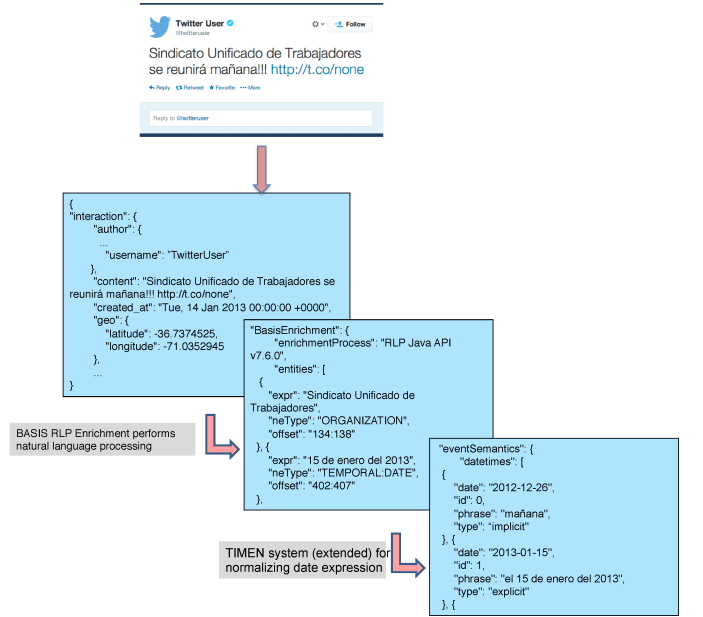


Figure 1: Message Enrichment

Applying BASIS technologies' Rosette Language Processing (RLP) tools, the language of the text is identified, the natural language content is tokenized and lematized and the named entities identified and classified. Date expressions are normalized and deindexed (using the TIMEN [37] package). Finally, messages are geocoded with a specification of the location (city, state, country), being talked about in the message. An example of this enrichment processing can be seen in Fig. 1.

We make use of different geocoding methodologies for geocoding news/blogs and twitter.

5.1 Geo-Coding

5.1.1 News/Blogs

To extract the protest location from news articles, we use *probabilistic soft logic* (PSL) [19] to build a model that performs robust, probabilistic inference given noisy signals.

PSL takes a set of weighted, logic-like rules and converts them into a continuous probability distribution over the unknown truth values of logical facts. These truth values in PSL are relaxed into the $[0, 1]$ interval. We use this mechanism to build a model that infers the semantic location of an article by weighing evidence coming from the Basis entity extractions and information in the World Gazetteer.

The primary rules in the model encode the effect that Basis-extracted location strings that match to gazetteer aliases are indicators of the article’s location, whether they be country, state, or city aliases. Each of these implications is conjuncted with an prior for ambiguous, overloaded aliases that is proportional to the population of the gazetteer location. For example, if the string “Los Angeles” appears in the article, it could refer to either Los Angeles, California, or Los Ángeles in Argentina or Chile. Given no other information, our model would infer a higher truth value for the article referring to Los Angeles, California, because it has a much higher population than the other options.

$$\begin{aligned} ENTITY(L, location) \tilde{\wedge} REFERSTO(L, locID) \\ \rightarrow PSLOCATION(Article, locID) \end{aligned}$$

$$\begin{aligned} ENTITY(C, location) \tilde{\wedge} IsCountry(C) \\ \rightarrow ArticleCountry(Article, C) \end{aligned}$$

$$\begin{aligned} ENTITY(S, location) \tilde{\wedge} IsState(S) \\ \rightarrow ArticleCountry(Article, S) \end{aligned}$$

The secondary rules, which are given half the weight of the primary rules, perform the same mapping of extracted strings to gazetteer aliases, but for extracted persons and organizations. Strings describing persons and organizations often include location clues (e.g., “mayor of Buenos Aires”), but intuition suggests the correlation between the article’s location and these clues may be lower than with location strings.

$$\begin{aligned} ENTITY(O, organization) \tilde{\wedge} REFERSTO(O, locID) \\ \rightarrow PSLOCATION(Article, locID) \end{aligned}$$

$$\begin{aligned} ENTITY(O, organization) \tilde{\wedge} IsCountry(O) \\ \rightarrow ArticleCountry(Article, O) \end{aligned}$$

$$\begin{aligned} ENTITY(O, organization) \tilde{\wedge} IsState(O) \\ \rightarrow ArticleCountry(Article, O) \end{aligned}$$

Finally, the model includes rules and constraints to require consistency between the different levels of geolocation, making the model place higher probability on states with its city contained in its state, which is contained in its country. As a post-processing step, we enforce this consistency explicitly by using the inferred city and its enclosing state and country, but adding these rules into the model makes the probabilistic inference prefer consistent predictions, enabling it to combine evidence at all levels.

$$\begin{aligned} PSLOCATION(Article, locID) \tilde{\wedge} Country(locID, C) \\ \rightarrow ArticleCountry(Article, C) \end{aligned}$$

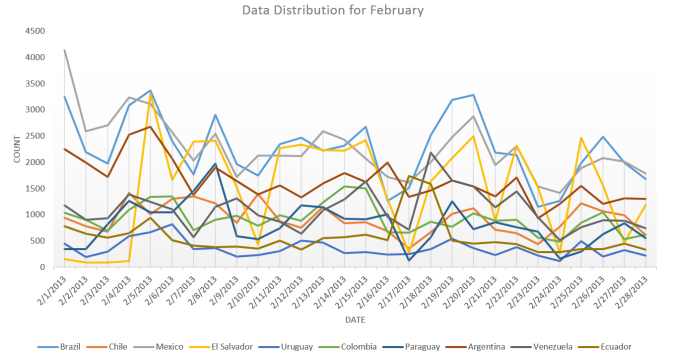


Figure 3: Rate of Arrival of News/Blogs

$$\begin{aligned} PSLOCATION(Article, locID) \tilde{\wedge} Admin1(locID, S) \\ \rightarrow ArticleState(Article, S) \end{aligned}$$

5.1.2 Twitter

The Twitter[2] geocoding is achieved by first considering the most reliable but least available source, viz. geotags, which give us exact geographic locations that can be reverse geocoded into place names. Second, we consider Twitter places and use place names present in these fields to geocode the place names into geographical coordinates. Finally, we consider the text fields contained in the user profile (location, description) as well as the tweet text itself to find mentions of relevant locations which can then be geocoded into geographical coordinates.

5.1.3 Facebook

We make use of only the facebook event data for our experiments. Facebook events that have a venue are only used. A venue of a facebook event generally contains a latitude, longitude, country, state, city, street, etc. Under cases where only latitude and longitude is given we do a reverse-geocoding by a KD-Tree lookup[1] from the World Gazetteer to get the country, state, city information.

6. APPROACH

All News, Blogs and Tweets are first searched for the phrases learned. Then the filtered documents are searched for the presence of a reference to a Future Date. In case of News/blogs we search for the Presence of a reference to a Future Date only within the sentence where the phrase was found to reduce error. For tweets, we search the entire tweet for the reference of a future date.

Then, finally, a warning/alert is issued for those documents which contains a location information. In the case of tweets, we found that issuing an alert from just one tweet lead to a lot of wrong alerts. We thus, further filter the tweets by setting a threshold (set to 5) on the number of re-tweets of the tweet under consideration.

Each Individual step is discussed in detail below.

6.1 Key Phrase Extraction

In Key Phrase Extraction, an input document is searched for the presence of one or more key phrases obtained in semi-automatic approach as given in 6.1.1. Each key phrase

WESTON, Fla — In December 2002, Ariel Dunaevski, then the owner of a furniture business in Caracas, Venezuela was on vacation in New York with his family when opponents of President Hugo Chávez called a crippling labor strike hoping to bring the government to its knees. As the protest wore on, paralyzing the country's oil industry and devastating the economy, the Dunaevskis saw a very uncertain future for Venezuela and arrived at a painful decision: they would be better off staying in the United States.

They flew to Florida and rented a house here in Weston, a suburb west of Fort Lauderdale that has become so popular with Venezuelan immigrants, it is known as Westonzuela.....

Venezuelans are outnumbered in South Florida by Cubans, Puerto Ricans, Colombians, Mexicans, Nicaraguans and Dominicans, according to data from the 2006 census, but Venezuelan leaders here believe their population may have vaulted to fourth place on that list, upwards of 100,000, taking into account those who have overstayed tourist visas.

"For a while you may forget about Chávez, forget about Miami, you're drinking your beer, you're insulting everybody, you're having fun," he said. "It's a way to forget about everything."

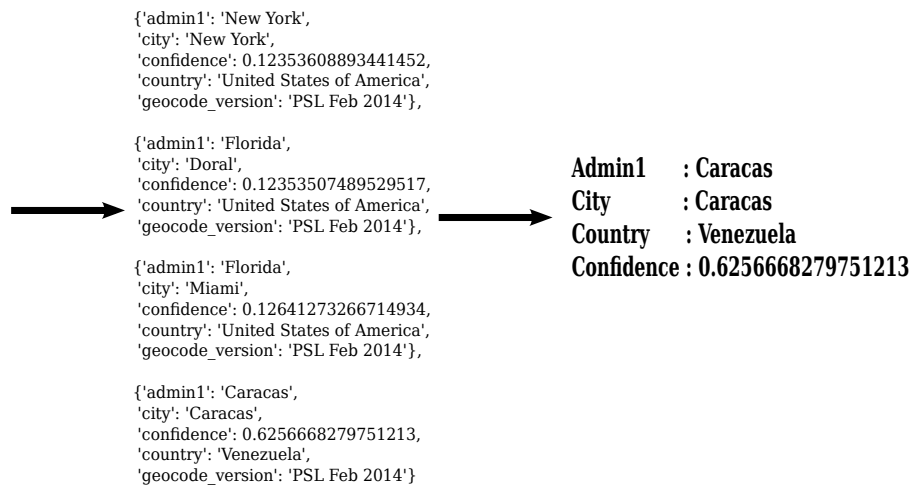


Figure 2: Red circles denote named entities identified as locations and blue denotes other types of entities. The article is reported from Weston Florida US and talks about the recent increase of venezuelan population in the US compared to other Latin American Nations like Cuba etc.[Sathappan writes: *TODO: Replace with a better protest example*]

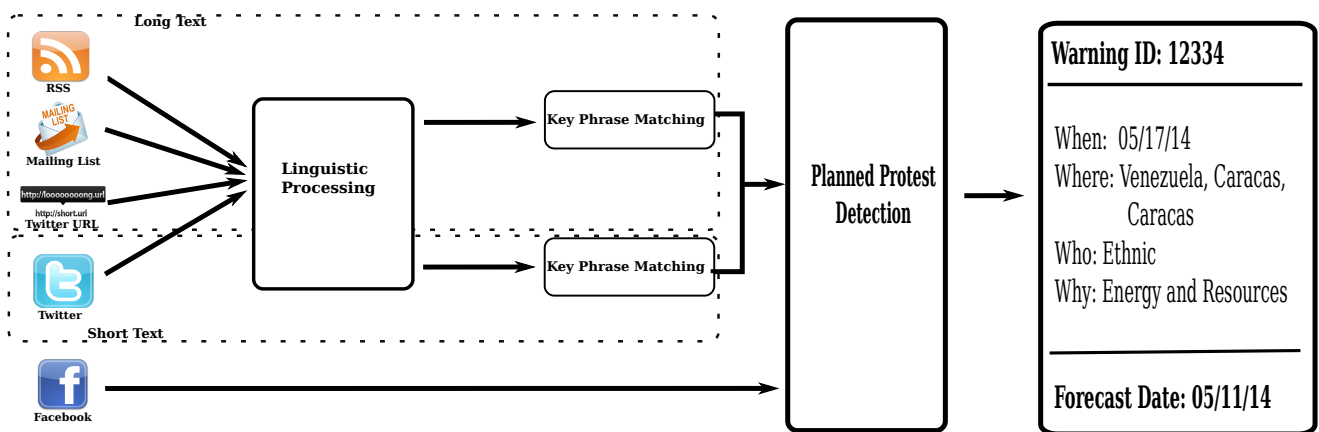


Figure 4: A diagram showing various steps of the Model

mostly consists of atleast two words. For checking the presence of a keyphrase, we search each individual word of the keyphrase independently in the input document with the constraint that they have to be within the same sentence and separated by atmost n words. n is the average distance between the keyphrase words in the learning phase.

6.1.1 phraselearning

We make use of a semi-automated approach to identify key-phrases for the purpose of extracting events from our data feeds. We use different sets of phrases for News and Twitter. We found that key-phrases provide a more accurate extraction than individual keywords.

Initially, a few seed phrases were obtained manually with the help of subject matter experts.

[Sathappan writes: *Jaime's Text* →]

The idea was to have a simple query list for future civil unrest actions, not trying to identify the type of protest. The analysis of news reports for planned protests on the printed media helped to create a minimum set of words to use in the query. We choose four nouns from the basic query that would indicate a civil unrest in all Latin American countries - *Demonstration, March, Protest and Strike*. We translated them into Spanish and Portuguese, including synonyms. We then set some verbs that indicated a future action and included the proper future conjugations - to organize, to prepare, to plan, to announce, etc. For twitter, shorter terms were used to overcome its limitations of characters. These phrases had a more direct call for action, different from the ones chosen for RSS feeds. For example - *Marchar, manhã de mobilização, vamos protestar, Huelga* - was identified.

These phrases were then parsed using a dependency parser and the grammatical relationship between the core subject word—*protest, manifestación, Huelga*, etc.—and any accompanying word - *plan, call, anunciar* - was extracted. These grammatical relations serve as extraction patterns as in [43]. Then, we build a corpora by finding out sentences that contain both a future date and the protest or its synonyms in any of the three languages. These sentences are then used as candidates for phrase list expansion. The phrase list is populated by finding phrases that follows the extraction pattern. The phrase learning is shown in Fig. /reffig:phraselearning

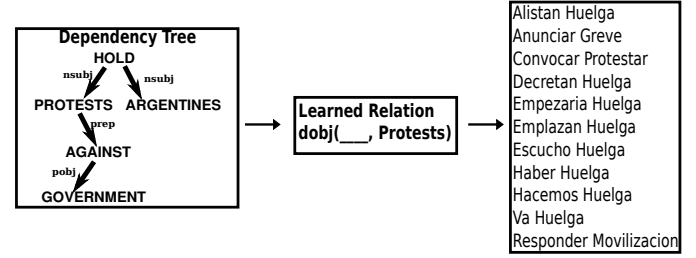
This final set of phrases, is then cleansed by an expert to get the final set of phrases. Around 122 phrases were obtained for News/Blogs and

To extend the initial set of phrases, a set of sentences/tweets containing a subject word and a future time/date expression was collected and parsed. This set of sentences was used to expand the set of planned protest phrases by extracting all keyword combinations that have the same grammatical relation with respect to the core subject word. The final set of planned protest phrases is then obtained after a manual revision of the phrases obtained in the last step.

By this approach, we learned 112 phrases for News/blogs and 156 for tweets.

The learned phrases are then used to filter the incoming stream of Documents (news/blogs/twitter). The phrases matching is done by first splitting the incoming document into sentences and then looking for the presence of each individual word of a key-phrase (by lemma) separated by a pre-fixed maximum offset-distance (set to 3). This methodology greatly increases processing speed.

Figure 5: An Example of Phrase Learning



6.2 classification

For News/Blogs and Facebook, we make use of Text Based Naive Bayes Classifier to identify the event-type and population. Unigram and Bi-gram word features are used for training the classifier.

For Twitter, as we send alerts based on a single tweet, we chose the event-type and population based on prior likelihood for that location.

7. EXPERIMENTS

We evaluate the model under a strict and a relaxed condition. Under strict evaluation an alert can be matched to an event only if there is city level location match and the forecasted event date is same as the true event date. In the relaxed evaluation, we allow the matching to happen if the alert and the GSR event are within a 300 KM radius and the forecasted event date lies within a given interval of the true event date. We try different matching intervals ranging from 0 to 7. If x is the allowed interval, then a matching can happen if an alert's forecasted event date is within $+x$ (upper bound) days of the GSR event and with a lower bound of $\min(x, GSREventDate - DateofForecast)$:

• Performance over the months

Fig. 7d provides the evaluation results of the model over the months with a source level breakdown. The QS reported is the weighted average of QS of all 10 countries where the weight for a country is the number of GSR events for that country. Twitter has a higher QS as multiple re-tweets of mention of future events in twitter is a direct indicator of the popularity of an event as well as the intent of people to join an event. While mention of Future events in News is simply a reporting of the event not much can be understood about the popularity of the event or about the people's support for the event.

• Country-wise performance

Table 3 presents the performance of the planned protest model (for March 2014) for each of the 10 countries of interest. It also presents a source wise breakdown. From the table it is evident that different data source prefer different countries. The News/Blogs data source produces alerts for most countries and also provides much higher recall as opposed to twitter which provides High QS but very little recall. Also News/Blogs has a higher lead-time of 4.57 days as compared to twitters 2.82.

• Case Study: Venezuelan and Brazilian Protests

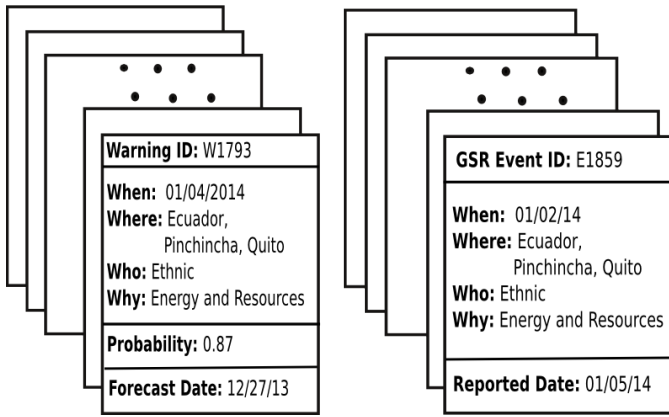


Figure 8: Structure of an Alert

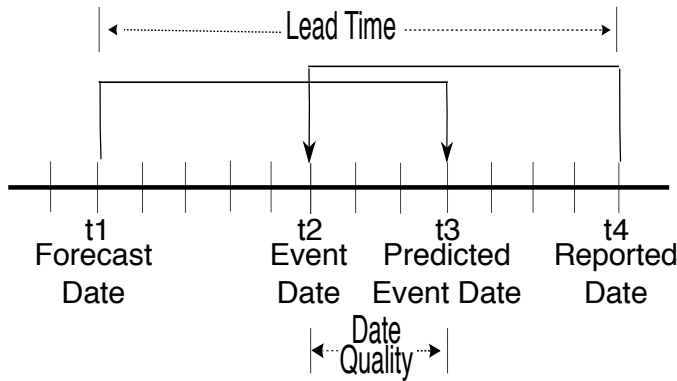


Figure 9: Matching Timeline

The recent Venezuelan protests against President Nicolas Maduro and the Brazilian Protests during June 2013 against Bus Fare Hike were two significant protests during our period of evaluation. Fig. 7b and Fig. 7a show how well the planned protest model was able to predict the unfolding of events under both situations. Fig. 7c showcases the ability of the model to forecast the violent events also.

• Lead-Time vs Quality Trade-Off

Fig. 7e shows that the QS of the planned protest model increases with time.

• Performance under stringent matching Criteria

Fig. 7g shows the performance of the model when the matching window is varied from 7 to 1 in steps. We can see that the model is not affected badly even under the strict matching interval of 1-day difference.

• Quality Score Distribution

The hump on the right side of the Fig. 7f signifies that a majority of the planned protest alerts are high quality.

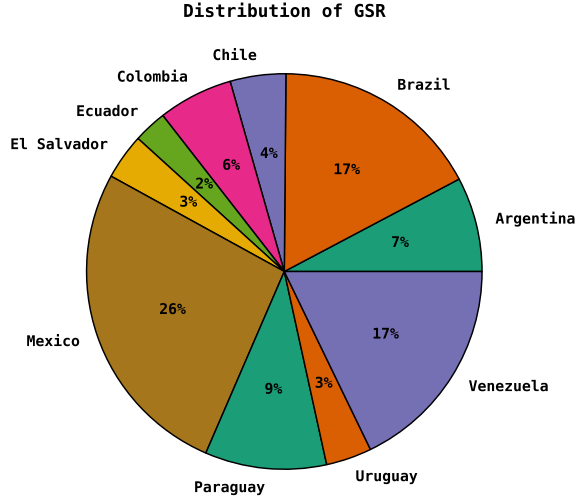
8. REFERENCES

- [1] <http://docs.scipy.org/doc/scipy-0.13.0/reference/generated/scipy.spatial.kdtree.html>.
- [2] <http://twitter.com>.

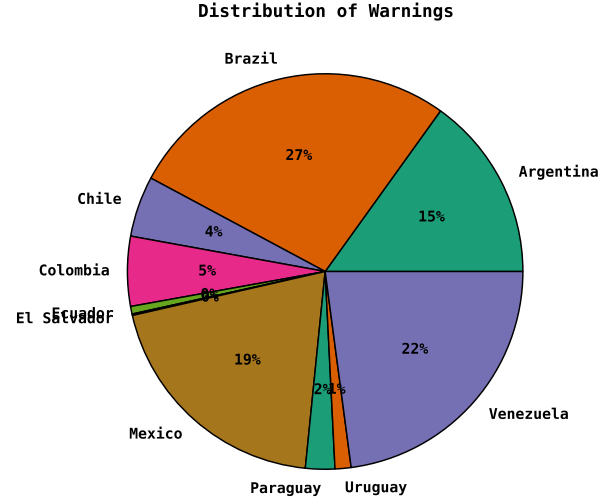
- [3] J. Allan, editor. *Topic Detection and Tracking: Event-based Information Organization*. Kluwer Academic Publishers, Norwell, MA, USA, 2002.
- [4] E. Amitay, N. Har'El, R. Sivan, and A. Soffer. Web-a-where: geotagging web content. In *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 273–280. ACM, 2004.
- [5] S. Bach, M. Broecheler, L. Getoor, and D. O'leary. Scaling mpe inference for constrained continuous markov random fields with consensus optimization. In *Advances in Neural Information Processing Systems*, pages 2663–2671, 2012.
- [6] S. H. Bach, M. Broecheler, S. Kok, and L. Getoor. Decision-driven models with probabilistic soft logic. 2010.
- [7] R. Baeza-Yates. Searching the future. In *SIGIR Workshop MF/IR*, 2005.
- [8] M. Banko, M. J. Cafarella, S. Soderland, M. Broadhead, and O. Etzioni. Open information extraction from the web. In *IN IJCAI*, pages 2670–2676, 2007.
- [9] H. Becker, F. Chen, D. Iter, M. Naaman, and L. Gravano. Automatic identification and presentation of twitter content for planned events. In *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media (ICWSM)*, page 2011.
- [10] H. Becker, D. Iter, M. Naaman, and L. Gravano. Identifying content for planned events across social media sites. In *Proceedings of the Fifth ACM International Conference on Web Search and Data Mining, WSDM '12*, pages 533–542, New York, NY, USA, 2012. ACM.
- [11] H. Becker, D. Iter, M. Naaman, and L. Gravano. Identifying content for planned events across social media sites. In *Proceedings of the fifth ACM international conference on Web search and data mining*, pages 533–542. ACM, 2012.
- [12] H. Becker, D. Iter, M. Naaman, and L. Gravano. Identifying content for planned events across social media sites. In *WSDM 2012*. ACM Press, 2012.
- [13] J. M. Box-Steffensmeier and B. S. Jones. *Event History Modeling: A Guide for Social Scientists*. 2004.
- [14] J. Braams. Babel, a multilingual style-option system for use with latex's standard document styles. *TUGboat*, 12(2):291–301, June 1991.
- [15] D. Braha. Global civil unrest: contagion, self-organization, and prediction. *PLoS One*, 7(10):e48596, Jan. 2012.
- [16] P. T. Brandt, J. R. Freeman, and P. A. Schrodt. Real time, time series forecasting of inter-and intra-state political conflict. *Conflict Management and Peace Science*, 28(1):41–64, 2011.
- [17] M. Brocheler, L. Mihalkova, and L. Getoor. Probabilistic similarity logic. *arXiv preprint arXiv:1203.3469*, 2012.
- [18] M. Broecheler and L. Getoor. Computing marginal distributions over continuous markov networks for statistical relational learning. In *Advances in Neural Information Processing Systems*, pages 316–324, 2010.
- [19] M. Broecheler, L. Mihalkova, and L. Getoor.

Table 3: Comparing forecasting accuracy of RSS vs Twitter

	News/Blogs				Twitter				Facebook				Combined			
	QS	Pr.	Rec.	LT	QS	Pr.	Rec.	LT	QS	Pr.	Rec.	LT	QS	Pr.	Rec.	LT
AR	3.14	0.32	0.69	3.94	3.52	0.78	0.14	3.14	3.70	0.50	0.04	3.00	3.02	0.36	0.80	4.50
BR	3.14	0.48	0.54	5.85	0.00	0.00	0.00	0.00	3.62	0.76	0.18	2.46	3.28	0.49	0.65	5.15
CL	3.06	0.91	0.67	5.40	3.52	1.00	0.23	4.29	0.00	0.00	0.00	0.00	3.16	0.83	0.80	5.92
CO	2.74	0.90	0.56	7.44	3.30	1.00	0.15	2.43	4.00	1.00	0.02	2.00	2.88	0.84	0.65	6.47
EC	0.00	0.00	0.00	0.00	2.32	1.00	0.06	17.00	0.00	0.00	0.00	0.00	2.32	0.50	0.06	17.00
MX	2.96	0.88	0.25	3.69	3.14	1.00	0.02	1.43	3.72	0.67	0.01	2.00	3.00	0.87	0.27	3.51
SV	3.22	1.00	0.03	1.0	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	3.22	1.0	0.03	1.0
PY	3.38	1.00	0.16	9.11	3.84	1.00	0.04	11.40	3.96	1.00	0.01	2.00	3.60	0.96	0.20	9.35
UY	3.24	1.00	0.29	2.40	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	3.24	1.00	0.29	3.24
VE	3.80	1.00	0.36	3.27	3.68	0.97	0.33	2.39	0.00	0.00	0.00	0.00	3.64	0.99	0.69	2.88
ALL	3.34	0.69	0.35	4.57	3.62	0.97	0.15	2.82	3.66	0.74	0.03	2.44	3.36	0.73	0.51	4.08



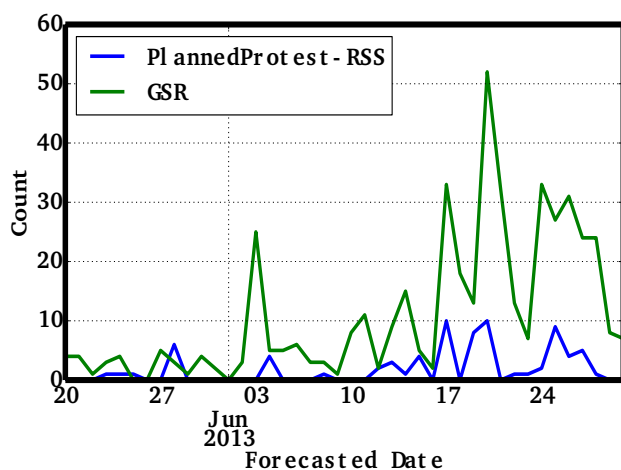
(a) GSR Distribution From 2012-11 to 2014-03



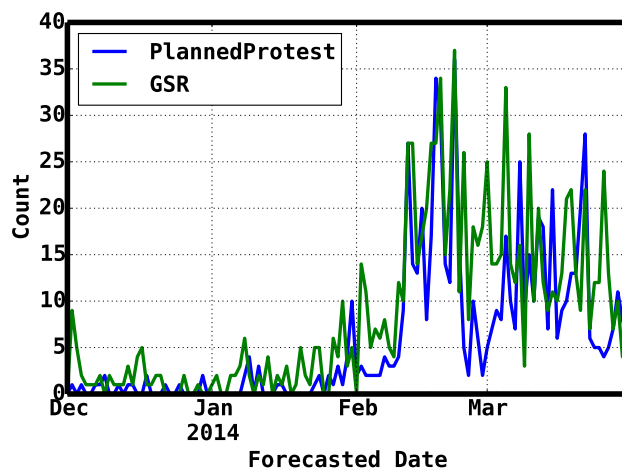
(b) Alerts Distribution From 2012-11 to 2014-03

Figure 6: Distribution of Alerts and GSR

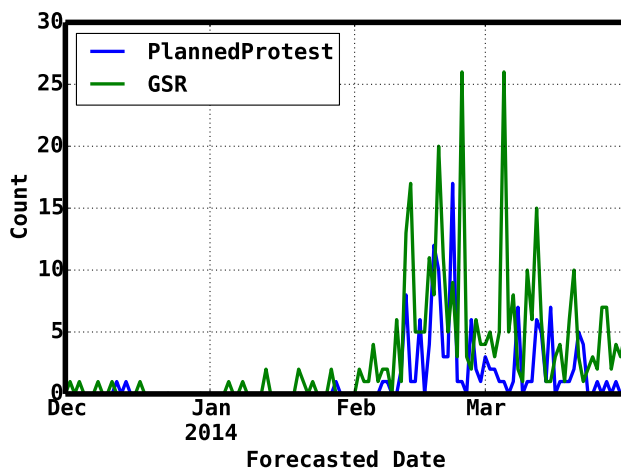
- Probabilistic similarity logic. In *Conference on Uncertainty in Artificial Intelligence*, 2010.
- [20] N. Chambers and D. Jurafsky. Template-based information extraction without the templates. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1, HLT '11*, pages 976–986, Stroudsburg, PA, USA, 2011. Association for Computational Linguistics.
- [21] R. Compton, C. Lee, T.-C. Lu, L. De Silva, and M. Macy. Detecting future social unrest in unprocessed twitter data: “emerging phenomena and big data”. In *Intelligence and Security Informatics (ISI), 2013 IEEE International Conference on*, pages 56–60. IEEE, 2013.
- [22] G. Dias, R. Campos, and A. Jorge. Future retrieval: What does the future talk about. In *Proceedings of Workshop on Enriching Information Retrieval of the 34th ACM Annual SIGIR Conference, SIGIR*, 2011.
- [23] G. Flores-Macías. Mexico’s 2012 elections: The return of the pri. *Journal of Democracy*, 24(1):128–141, 2013.
- [24] E. Gabrilovich, S. Dumais, and E. Horvitz. Newsjunker: Providing personalized newsfeeds via analysis of information novelty. In *Proceedings of the 13th International Conference on World Wide Web, WWW '04*, pages 482–490, New York, NY, USA, 2004. ACM.
- [25] S. González-Bailón, J. Borge-Holthoefer, A. Rivero, and Y. Moreno. The dynamics of protest recruitment through an online network. *Sci. Rep.*, 1:197, Jan. 2011.
- [26] P. Howard, A. Duffy, D. Freelon, M. Hussain, W. Mari, and M. Marzaid. Opening Closed Regimes: What Was the Role of Social Media During the Arab Spring? 2011.
- [27] T. Hua, C.-T. Lu, N. Ramakrishnan, F. Chen, J. Arredondo, D. Mares, and K. Summers. Analyzing Civil Unrest through Social Media. *Computer*, 46(12):80–84, Dec. 2013.
- [28] B. Huang, A. Kimmig, L. Getoor, and J. Golbeck. Probabilistic soft logic for trust analysis in social



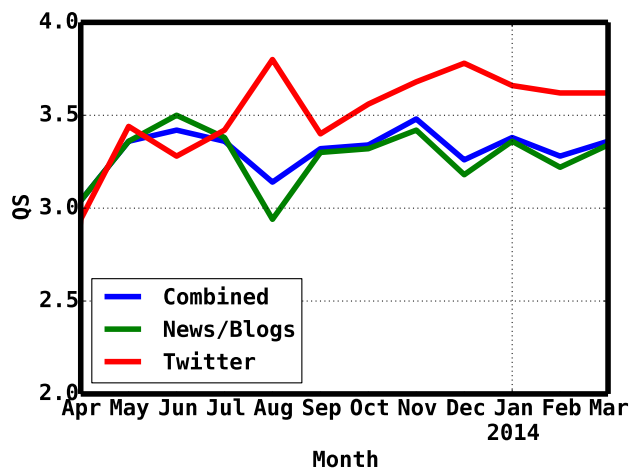
(a) System Performance during Brazilian Spring



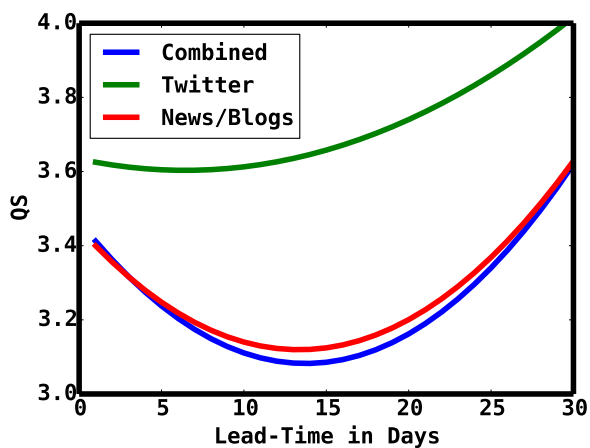
(b) Venezuelan Protests



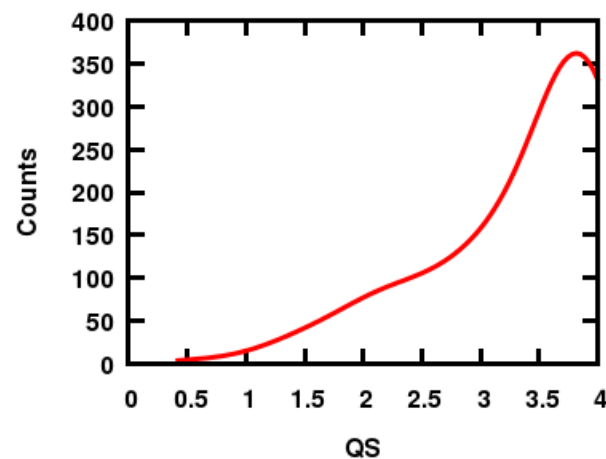
(c) Venezuelan Violent Protests



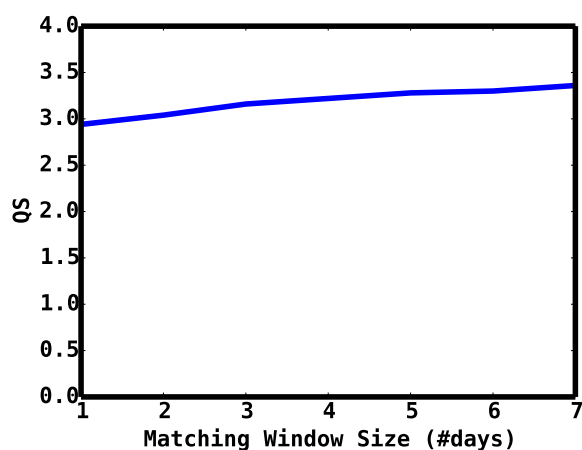
(d) Quality Score over the months



(e) Lead-Time vs Quality Score



(f) Quality Score Distribution



(g) QS vs Matching Interval Trade-Off

- networks. In *International Workshop on Statistical Relational AI*, 2012.
- [29] A. Jatowt and C.-m. Au Yeung. Extracting collective expectations about the future from large text collections. In *Proceedings of the 20th ACM International Conference on Information and Knowledge Management*, CIKM '11, pages 1259–1264, New York, NY, USA, 2011. ACM.
 - [30] N. Kallus. Predicting crowd behavior with big public data. *CoRR*, abs/1402.2308, 2014.
 - [31] N. Kallus. Predicting Crowd Behavior with Big Public Data. Feb. 2014.
 - [32] H. Kawai, A. Jatowt, K. Tanaka, K. Kunieda, and K. Yamada. Chronoseeker: Search engine for future and past events. In *Proceedings of the 4th International Conference on Ubiquitous Information Management and Communication*, ICUIMC '10, pages 25:1–25:10, New York, NY, USA, 2010. ACM.
 - [33] A. Kimmig, S. Bach, M. Broecheler, B. Huang, and L. Getoor. A short introduction to probabilistic soft logic. In *Proceedings of the NIPS Workshop on Probabilistic Programming: Foundations and Applications*, pages 1–4, 2012.
 - [34] K. Leetaru and P. Schrodtt. GDELT: Global data on events, location, and tone, 1979–2012. *ISA Annual Convention*, pages 1979–2012, 2013.
 - [35] K. H. Leetaru. Culturomics 2.0: Forecasting large-scale human behavior using global news media tone in time and space. *First Monday*, 16:1–22, 2011.
 - [36] 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*, SIGIR '05, pages 106–113, New York, NY, USA, 2005. ACM.
 - [37] H. Llorens, L. Derczynski, R. J. Gaizauskas, and E. Saquete. TIMEN: An open temporal expression normalisation resource. In *LREC*, pages 3044–3051. European Language Resources Association, 2012.
 - [38] A. Memory, A. Kimmig, S. Bach, L. Raschid, and L. Getoor. Graph summarization in annotated data using probabilistic soft logic. *status: accepted*, 2012.
 - [39] S. P. O'Brien. Crisis Early Warning and Decision Support: Contemporary Approaches and Thoughts on Future Research. *Int. Stud. Rev.*, 12(1):87–104, Mar. 2010.
 - [40] R. Pape. *Dying to Win: The Strategic Logic of Suicide Terrorism*. Random House Trade Paperbacks, 2006.
 - [41] K. Radinsky and E. Horvitz. Mining the web to predict future events. In *Proceedings of the Sixth ACM International Conference on Web Search and Data Mining*, WSDM '13, pages 255–264, New York, NY, USA, 2013. ACM.
 - [42] E. Riloff and J. Wiebe. Learning extraction patterns for subjective expressions. In *Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing*, EMNLP '03, pages 105–112, Stroudsburg, PA, USA, 2003. Association for Computational Linguistics.
 - [43] E. Riloff and J. Wiebe. Learning extraction patterns for subjective expressions. In *Proceedings of the 2003 conference on Empirical methods in natural language processing*, pages 105–112. Association for Computational Linguistics, 2003.
 - [44] A. Ritter, O. Etzioni, and S. Clark. Open domain event extraction from twitter. In *ACM SIGKDD 2012*, 2012.
 - [45] A. Ritter, Mausam, O. Etzioni, and S. Clark. Open domain event extraction from twitter. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '12, pages 1104–1112, New York, NY, USA, 2012. ACM.
 - [46] 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*, WWW '10, pages 851–860, New York, NY, USA, 2010. ACM.
 - [47] J. Sankaranarayanan, H. Samet, B. E. Teitler, M. D. Lieberman, and J. Sperling. Twitterstand: News in tweets. In *Proceedings of the 17th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, GIS '09, pages 42–51, New York, NY, USA, 2009. ACM.
 - [48] P. A. Schrodtt. Automated production of high-volume, near-real-time political event data. In *American Political Science Association meetings*, 2010.
 - [49] P. D. Senese and J. A. Vasquez. *The Steps to War: An Empirical Study*. Princeton University Press, 2008.
 - [50] J. D. Singer, editor. *The Correlates of War I: Research Origins and Rationale*. Free Press, 1979.
 - [51] H. Tops, A. van den Bosch, and F. Kunneman. Predicting time-to-event from twitter messages. 2013.
 - [52] S. Truve. Big data for the future: Unlocking the predictive power of the web(white paper).
 - [53] S. Truve. Big data for the future: Unlocking the predictive power of the web(white paper).
 - [54] A. H. urriyetoglu, F. Kunneman, and A. van den Bosch. Estimating the time between twitter messages and future events. pages 20–23, 2013.
 - [55] T. Y. Wang. *Arms Transfers and Coups d'Etat: A Study on Sub-Saharan Africa*, 1998.
 - [56] J. Xu, T.-C. Lu, R. Compton, and D. Allen. Civil unrest prediction: A tumblr-based exploration. In *Social Computing, Behavioral-Cultural Modeling and Prediction*, pages 403–411. Springer, 2014.
 - [57] Y. Yang, T. Pierce, and J. Carbonell. A study of retrospective and on-line event detection. In *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '98, pages 28–36, New York, NY, USA, 1998. ACM.