

Predicting Civil Unrest by Categorizing Dutch Twitter Events

Rik van Noord¹✉, Florian A. Kunneman², and Antal van den Bosch^{2,3}

¹ CLCG, University of Groningen, Groningen, The Netherlands
r.i.k.van.noord@rug.nl

² Centre for Language Studies, Radboud University, Nijmegen, The Netherlands
{f.kunneman,a.vandenbosch}@let.ru.nl

³ Meertens Institute, Amsterdam, The Netherlands

Abstract. We propose a system that assigns topical labels to automatically detected events in the Twitter stream. The automatic detection and labeling of events in social media streams is challenging due to the large number and variety of messages that are posted. The early detection of future social events, specifically those associated with civil unrest, has a wide applicability in areas such as security, e-governance, and journalism. We used machine learning algorithms and encoded the social media data using a wide range of features. Experiments show a high-precision (but low-recall) performance in the first step. We designed a second step that exploits classification probabilities, boosting the recall of our category of interest, social action events.

Keywords: Civil unrest · Event categorization · Event detection

1 Introduction

Many instabilities across the world develop into civil unrest. Unrest often materializes into crowd actions such as mass demonstrations and protests. A prime example of a mass crowd action in the Netherlands was the *Project X* party in Haren, Groningen, on September 21, 2012. A public Facebook invitation to a birthday party of a 16-year old girl ultimately led to thousands of people rioting [18]. The riots could only be stopped by severe police intervention, resulting in more than 30 injuries and up to 80 arrests. Afterwards it was concluded that the police were insufficiently prepared and that they were not well enough informed about the developments on social media. An evaluation committee recommended the development of a nation-wide system able to analyze and detect these threats in advance [13]. In this paper, we describe a system that leverages posts on Twitter to automatically predict such civil unrest events before they happen.

To facilitate this objective we start from a large set of open-domain events that were automatically detected from Twitter from a period spanning multiple

years by the approach described in [9]¹. From this set we aim to identify the events that are materializations of civil unrest, henceforth *social action events*. Arguably, a system that detects social actions should not only reliably detect events where large groups of people come together, but should also exclude events where people gather for different reasons (e.g. soccer matches, music performances) and for which authorities are sufficiently prepared. In addition, a system able to detect several categories reliably might be useful for other applications as well, such as presenting tourists with events of a certain type and in a certain time range. For these reasons, instead of focusing on this event type only, we categorize all events into a broad categorization of events, and distinguish social actions as one of the event types.

This paper is structured as follows. In Sect. 2 we provide an overview of related work, discussing both the tasks of predicting social action events and categorizing Twitter events. In Sect. 3 we present the experimental set-up, discussing the data, event annotations and event classification. We present the results of our system evaluation in Sect. 4, and analyze the retrieval of social action events, as well as the most informative features, in Sect. 5. Conclusions and a discussion are presented in Sect. 6.

2 Related Work

2.1 Predicting Social Action Events

There is a small body of work dedicated to detecting social action events. [3] aim to predict civil unrest in South America based on Twitter messages. In contrast to our approach, they predict such events directly from tweets, by matching them with specific civil unrest related keywords, a date mention, and one of the predefined locations of interest. Their system obtains a precision of 0.55 on a set of 283 predefined events. The main drawback of their approach is that it has no predictive abilities. For example, the system is not able to detect social action events that use newly emerging keywords for a specific event, or take place in a new location. As a consequence, their system has a low recall; many future social actions are likely to go undetected.

A more generic approach to detecting social action events is the EMBERS system by [16]. They try to forecast civil unrest by using a number of open source data sources such as Facebook, Twitter, blogs, news media, economic indicators, and even counts of requests to the TOR browser.² Using multiple models, the system issues a warning alert when it believes a social action event is imminent. Tested over a month, the EMBERS system attained a precision of 0.69 and a recall of 0.82. [6] provide a more detailed explanation of some of the EMBERS models, reporting similar F-scores as [16]. They also compare

¹ A live event detection system using the method of [9] is available at <http://lamaevents.cls.ru.nl/>.

² TOR requests are an indication of the number of people who choose to hide their identity and location.

the impact of different data sources, concluding that Social Media (including Twitter, blogs and news) is the most informative source. [12] test EMBERS when only taking Twitter information into account, reporting a precision of 0.97 but a recall of 0.15.

2.2 Categorizing Events

Some approaches based on Twitter perform some form of broad categorization of events [15,20]. These approaches either identify which topics are often talked about on Twitter, or focus on the categorization of users instead of events. To our knowledge, the only approach that focuses on the categorization of automatically detected events is the one proposed by [17]. They apply Latent Dirichlet Allocation [2] to a set of 65 million events to generate 100 topical labels automatically. Manual post-annotation winnowed these down to a set of 37 meaningful categories. 46.5% of the events belong to one of these categories, while 53.5% of the events are in a rest category. [17] compared their unsupervised approach to categorizing Twitter events to a supervised approach. They selected the best 500 events (detected with the highest confidence) and manually annotated them by event type. Their unsupervised approach obtained an F1-score of 0.67, outperforming the supervised approach which obtained an F1-score of 0.59. However, they do show that the F1-score of the supervised approach steadily increases when using more training instances.

3 Experimental Set-Up

Our study starts with a set of automatically detected events from Twitter, described in Sect.3.1. We manually annotate two subsets of these events by type, and subsequently train a machine learning classifier on several feature types extracted from these events. Performance is both evaluated on the annotated event sets and on the larger set of remaining events.

3.1 Data

Event Set. To perform automatic event categorization, we use the event set described in [8] which was extracted based on the approach described in [9]. As this approach was applied to Dutch tweets, the set mainly comprises Dutch events. This approach, based on the method of [17], comprises the extraction of explicit time expressions and entities from tweets, identifying date-entity pairs as event when they co-occur together for at least five times and display a good fit as measured by the $G2$ log likelihood ratio statistic.

An example of a detected social action event on Twitter is shown in Table 1 [19]. Each event has a set of attributes, such as the date, keywords, tweets and event score. The event score is linked to the size and popularity of an event. For the exact calculation of this score, we refer to [9, pp. 13]. Over a 6-year period (2010–2015), [8] ultimately obtained 93,901 events. This event set is used for our categorization system.

Table 1. An example of an actual Dutch social action event with five example tweets (translated to English).

Date	21-09-2013
Keywords	#demonstration, budget cuts
21 September: say no to the new budget cuts! #demonstration	
Are you also coming to the #demonstration #21september ? #action is necessary!	
Come Sept 21 to The Hague to demonstrate against the cabinet #demonstration	
It is allowed again tomorrow, no excuses not to go to #thehague #resistance	
8 days to go #PVV #demonstration against #cabinet at #koeplein The Hague!	

Event Annotations. We select two sets of events for manual labeling. Our first event set contains the 600 events with the highest event score in the output of [8]. This enables us to make an approximate comparison to [17], who evaluated their system on the basis of their 500 top-ranked events. We refer to the set of events with the highest event scores as the *best event set*.

Our second event set is created by randomly selecting an event from the ranked total event set for intervals of 155 events (with all events ranked by event score), excluding the best 600 events of the best event set. We refer to this event set as the *random event set*. Non-Dutch events were manually removed from both event sets, leaving 586 events in the best events set and 585 in the random events set.

First, we annotated the set of best events. Seven annotators were involved in the annotation process, who all at least annotated 40 and at most 175 of the events in the best event set. 195 of the 586 best events received a double annotation so that we are able to calculate inter-annotator agreement. The other 390 events, as well as the 585 random events, were annotated by one annotator. Similar to [17], the annotator is asked two questions for each event:

- Is this an actual event according to the definition?
- What is the category of this event?

We employ the same definition as [9] as to what constitutes an actual event: ‘An event is a significant thing that happens at some specific time and place’, where ‘significant’ is defined as ‘something that may be discussed in the news media’. An event in our full event set is not necessarily a proper event according to this definition, as the detection procedure makes errors. Since we are not interested in the category of a non-event, the events that are annotated as a non-event are filtered from the event set.

We defined ten possible categories after an initial manual inspection of about 200 events. They are listed in Table 2. *Social action* is the category of interest.

Table 2. The ten different categories with examples.

Category	Example events
Social action	Strikes, demonstrations, flashmobs
Sport	Soccer match, local gymnastics event
Politics	Election, public debate
Broadcast	Television show, premiere of a movie
Public event	Performance of a band, festival
Software	Release of game, release of new iPhone
Special day	Mother’s Day, Christmas
Celebrity news	Wedding or divorce of a celebrity
Advertisement	Special offers, retweet and win actions
Other	Rest category

As arguably less straightforward categories we included *special day* and *advertisement* because manual inspection of the data suggested that those types of events were frequent enough to deserve their own category.

The 195 events annotated by two coders yielded a Krippendorff’s alpha [5] of 0.81 on judging whether or not it was an actual event and 0.90 on categorizing events. These scores can be considered excellent [7] and show that we can reliably view the events that were annotated once as if they were annotated correctly. Therefore, the 586 random events could be annotated once by two annotators.

Events that were (at least once) annotated as a non-event are removed from the event set, as well as events where annotators disagreed on the category. 27.4% of the best events were a non-event, leaving 425 of the best events. In the random event set 38.1% were discarded as non-events, leaving 362 events.

The annotations by event category are shown in Fig. 1. *Public event* is the dominant category, comprising 29.6% of the best events and 44.5% of the random events. Most other event categories occur fairly regularly, except *advertisement* and *celebrity news*. The latter category was so infrequent that it was removed from both event sets. *Advertisement* was removed from the random event set, but was retained for the best event set.

3.2 Training and Testing

Based on the annotated events we trained a machine learning classifier to distinguish the ten event types. We describe the event features, classification approaches and evaluation below.

Feature Extraction. To enable the classifier to learn the specific properties of each event category we extract several types of features from each event. They are listed in Table 3. The first four feature types are derived from [9]. The first feature type, the event score, describes the link between the event keywords and

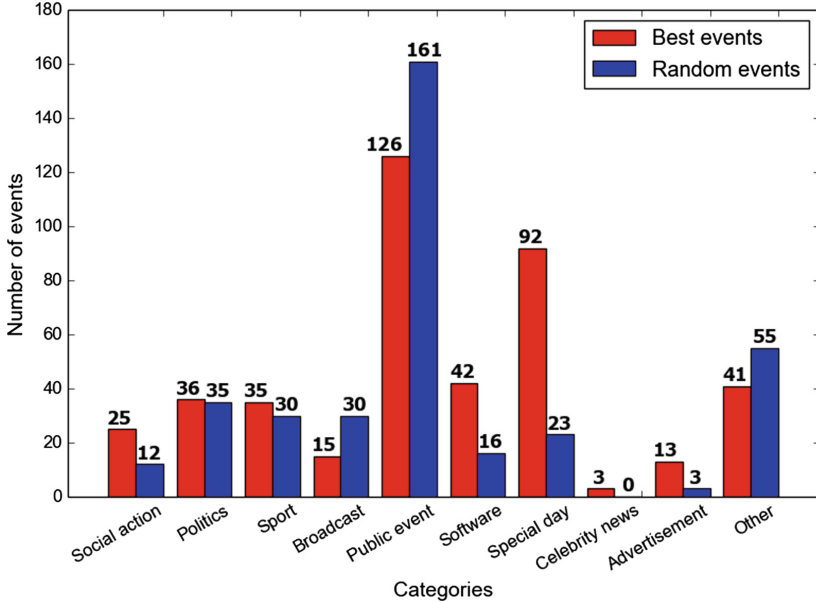


Fig. 1. The ten different categories with the number of annotated examples in the best and random event set.

the date of the event. This score gives an indication of the confidence that the set actually represents an event. Second, the keyword scores give an indication of the commonness of each event keyword, based on the commonness score as described in [10]. Third, the event date might help to recognize event types that are linked to big events, such as elections. Fourth, we extract the number of tweets, which might reflect the popularity of an event. We also distinguish between the numbers of tweets before, during, and after³ an event. Fifth, we extract each word used in the event tweets as a feature, jointly referred to as bag-of-words features. Such features provide the classifier with a lot of information, but there is no deeper reasoning involved concerning the words in question. The most informative words per category are shown in Table 4. As a sixth feature type we scored the average subjectivity and polarity of each event tweet, using the approach by [4]. The subjectivity and polarity score of the event are averaged over the scores of all event tweets. Some event types might be referred to fairly objectively in tweets, while others might stir more sentiment.

³ Since we wanted to provide our system with as much training data as possible, we also extracted relevant tweets that were posted after the event took place. Obviously, when predicting events in the future, this type of data will be unavailable.

Table 3. Types of extracted features with descriptions.

Feature type	Description
Event score	Single feature specifying the event-score
Keyword scores	Feature per keyword specifying the keyword-score
Event date	Single feature specifying the event date
Tweet count	Three features, specifying the total number of tweets and number of tweets before and after the event
Bag-of-words	Each unique word has its own feature, the value of each feature is determined by how often the word occurs in the tweets of the event
Sentiment	Two features: the average subjectivity and the average polarity, calculated over all tweets of the event
Periodicity	Two features: one binary feature that specifies if the event is periodic and one feature that specifies the periodicity type (e.g. yearly)
Wikipedia	Each unique Wikipedia type has its own feature, the value of each feature is determined by how often the type occurs in the tweets of the event

The seventh feature type indicates whether an event is of a periodic nature. This feature is based on the output from a periodicity detection system described by [8]. Finally, we employ DBpedia [1] in order to generalize over the different named entities present in the events. Since we want to generalize over the different terms, we are especially interested in the `type` attribute of the entity in DBpedia. This gives us a broader description of the named entities in question and therefore allows for generalization of previously unseen entities. For example, Feyenoord is a *SoccerClub*, *SportsTeam* and *Organisation*, while Justin Bieber is an *Artist*, *MusicalArtist*, *MusicGroup*, and a *NaturalPerson*. We extract the different DBpedia `types` for each event keyword. The keywords are linked to DBpedia using Wikification [11].

Table 4. An ordered list of the 8 most indicative words per category according to their $tf * idf$ score (translated from Dutch to English).

Category	Most indicative words
Social action	Against everyone protest respect they demonstration all
Politics	Votes elections vote cda vvd pvda d66 groenlinks
Sport	Match against soccer wins rt ajax psv tonight
Broadcast	Tv watch tonight episode see tvtip show season
Public event	rt tonight was today what who much tomorrow
Software	Apple iphone microsoft out gta wait comes windows
Special day	Today rt what everyone day on if celebrate
Other	Not rt that no what will so today one

Classification. Based on the extracted feature sets along with the annotated categories, we train a Naive Bayes classifier⁴ using the Python module Scikit-learn⁵ [14]. We use Laplace smoothing ($\alpha = 1.0$) and learn the prior probabilities per class. No correction method for document length is employed.

We applied two methods to increase the performance of the classifier: down-sampling the dominant *public event* class, and performing bag-of-words classification as a first-step classification. The first method simply reduces the number of *public events* in the training set to ensure it does not hinder the performance of the minority classes. The number of *public events* is reduced to the same frequency as the second-most frequent class in the event set, resulting in the deletion of 34 events in the best event set and 106 events in the random event set.

The second method only feeds the bag-of-words features to the classifier in an initial stage, and subsequently adds the resulting classification to the set of other features. The advantage of this stacking method is that it reduces the dominance of the word features compared to the other features, allowing the classifier to view the word features as a single source of information. Also, it enables us to measure the impact of the non bag-of-words features in comparison to a bag-of-words baseline.

Evaluation. The performance on categorizing events is evaluated in two ways. The first is to apply 5-fold cross validation on the annotated sets of events. We do this for both the best event set (for a comparison with [17]) and the random event set, calculating the average precision, recall and F1-score.⁶ The second way is to evaluate the results on a set that was never used in the training phase. The classifiers are trained on the two sets of annotated events and subsequently applied to the remaining 92,701 events. Performance on these unseen events is evaluated by manually inspecting a subset of them. As *public event* appeared to be a very dominant category, occurring 81,538 times in the full set of events according to the classifier, it was not feasible to randomly select a set of events to be used as evaluation set. This is why we focus on evaluating the precision of each classification category separately. We randomly selected 50 events per category for evaluation, except for our category of interest, social action events, for which we include all 93 events classified with this category. *Advertisement* could only be evaluated for 25 events, as it was only predicted 25 times. This ultimately resulted in a total set of 468 events, which we refer to as the *Evaluation set*.

⁴ In addition to Naive Bayes, we experimented with Support Vector Machines and K-nearest neighbors. We will only report on the outcomes of Naive Bayes, which yielded the best performance.

⁵ <http://scikit-learn.org>.

⁶ This was calculated by using the *weighted* setting in scikit-learn, which is why the F-score is not necessarily between precision and recall.

4 Results

4.1 Annotated Set

Table 5 shows the most important results of the 5-fold cross validation. Averaged over all categories, the best event set obtained an F1-score of 0.65, while the random event set received an F1-score of 0.58. It appears to be easier to classify events with a higher event score. However, we found no significant effect of event score when doing a least-squares logistic regression test for the random event set ($r(360) = -0.05$, $p = 0.39$). This suggests that there is a small subset of events with a very high event score that is easier to classify, but that there is no significant effect of event score in general.

Comparing the setting where only bag-of-words is used as a feature with the setting where the classification based on bag-of-words is added as a feature to the other features, the latter setting yields the best outcomes.

The score on our best event set is similar to the score of [17]. However, it is hard to make a fair comparison, since they did not include a category distribution of the test set in their 37-class problem.⁷

Social action is predicted at a high precision in the best events set, but the scores for the random events are poor. This might be due to the low number of instances in this set (12), in comparison with the 25 social actions in the best event set.

Table 5. The results of the 5-fold cross validation for the Naive Bayes algorithm while down-sampling the dominant *public event* class.

		All categories			Social actions		
		Prec	Rec	F1	Prec	Rec	F1
Best events	Only bag-of-words	0.67	0.59	0.55	0.68	0.41	0.52
	Bag-of-words as feature	0.67	0.67	0.65	0.79	0.44	0.56
Random events	Only bag-of-words	0.61	0.60	0.57	0.40	0.17	0.24
	Bag-of-words as feature	0.64	0.60	0.58	0.40	0.17	0.24

Down-sampling increased the F1-score by 0.05 for the best event set and 0.06 for the random event set.

4.2 Evaluation Set

Table 6 shows the results on the Evaluation set, listing the precision per category. In general, these scores are high for a 9-class classification task. The precision per class is even 1.00 for *sport* and *politics*, meaning that if the classifier predicted those categories, it did so perfectly. The categories *public event* and *advertisement* score below 0.70, however. The low precision for *public event* impacts the

⁷ In personal communication, we asked Alan Ritter about this distribution. Unfortunately, he was unable to recover the document with the specific division of categories in the test set.

Table 6. The precision and number of predicted instances per category.

Category	Instances	Precision	Category	Instances	Precision
Social action	93	0.80	Software	1,630	0.96
Politics	2,170	0.86	Special day	1,722	0.78
Sport	2,771	1.00	Advertisement	25	0.51
Broadcast	206	1.00	Other	1,535	0.70
Public event	81,538	0.57			

overall performance of the classification system substantially. As 81,538 out of 92,701 events were classified as a *public event*, a precision of 0.57 leads to about 35 thousand incorrectly classified events.

We should keep in mind that the non-events were not excluded from the full event set. It was estimated that 38.1% of all detected events are not events. In the training phase these non-events were excluded, so it is likely that the classifier will assign many non-events in the full event set to the most frequent category. A large part of the bias to *public event* may be due to the occurrence of non-events in the full event set. This leads us to conclude that if there were a more reliable way to automatically exclude non-events, the results of the general categorization would considerably improve.

The results for the *Social Action* category are promising, since the 93 *social actions* in this set were predicted with a precision of 0.80. However, we estimate that the recall of this category will be low. Only 93 out of 92,701 events (0.1%) were predicted as a *social action*, while 3.3% of events were annotated as a *social action* in the random event set.

5 Analysis

5.1 Increasing the Recall for Social Action Events

Our main goal is to detect *social action events* and possibly alerting the authorities when such an event will take place. Therefore, we rather show a large list of events that might be a social action event that actually includes most of the actual events, than a system that often misses them. Since we are not talking about thousands of events daily, an analyst could annotate the set of possible social action events manually. We thus prefer a high recall to a high precision. Therefore, we propose a method to increase the recall of social action events, at minimal precision costs.

In order to increase recall we make use of the Naive Bayes classifier probability by category that is assigned to each event. Events for which *social action* obtained the second highest probability are ignored by default, as another category is picked. One way to remedy this is to classify all events where *social action* was the second most probable class. We refer to these events as **secondary social action events**. By doing this we were able to expand this set with 226

additional events, which we annotated manually. 26 of the 226 *secondary social actions* were annotated as a non-event and were thus excluded from the set. 130 of the remaining 200 events were indeed annotated as a social action, resulting in a precision of **0.65**. Adding the 200 events to the *social action* events in the evaluation set results in a drop of total precision from 0.80 to 0.69. Thus, since we could add 130 *social action events* to the 56 that were already found as a primary classification, the recall was increased by **232%** while the precision only dropped by **14%**. Hence, including the *secondary social action events* seems a useful method for increasing the recall, while only mildly hurting precision.

A possible other method that exploits the actual Bayesian probabilities would be to select all events for which the probability of *social action event* exceeds a certain threshold. This would then allow us to pick a specific precision-recall trade-off, instead of simply relying on events that were second in the ranking of probabilities. Investigating this is left for future work.

5.2 Most Informative Features

In order to achieve some insight from the most informative features for the two event sets, we calculated the chi-squared value for each feature in relation to the category label. These are listed in Table 7. The most informative features are generally intuitive. They include words such as *stemmen (to vote)* and *stem (vote)* as indicators of a political event, but also specific hashtags such as *#VVD* and *#CDA*; CDA and VVD are political parties in The Netherlands. The best predictors for *sport* are the DBpedia type features *SoccerClub* and *ClubOrganization*. The most indicative features of the category *social action* are the words *protest* and *demonstratie (demonstration)*. Although these words almost exclusively occurred in *social action events*, due to their low frequency they do not rank in the feature top 100.

Table 7. The eight best features for the best and random event set, based on their chi-squared value. Non-word features are in italics. Features are only included if they occurred at least ten times in their event set.

Best events		Random events	
Feature	Category	Feature	Category
stemmen (vote)	Politics	<i>ClubOrganization</i>	Sport
stem (vote)	Politics	<i>SoccerClub</i>	Sport
<i>19-03-2014</i>	Politics	wint (wins)	Sport
<i>SoccerClub</i>	Sport	wedstrijd (match)	Sport
<i>#vvd</i>	Politics	2015	Politics
wedstrijd (match)	Sport	seizoen (season)	Broadcast
<i>ClubOrganization</i>	Sport	tv	Broadcast
<i>#cda</i>	Politics	tegen (against)	Sport

The polarity, subjectivity and periodicity features turned out to be less valuable, ranking in the bottom 25% of all features. This is surprising, since *special days* are often periodic, while it is, for example, uncommon for *social action events* to be periodic.

6 Conclusion and Discussion

In this study we presented a generic event categorization system which we evaluated particularly on its ability to predict civil unrest. The general categorization system has a bias towards the dominant category *public event*, but has a high precision for the other categories, including *social action*. The recall for *social action* was low; a follow-up step that exploited the specific per-class probabilities generated by the Naive Bayes classifier led to a considerable improvement in recall of 232%, at the minor cost of a 14% decrease in precision.

The study by [17] is the only related study that also produced an extensive evaluation of event categorization, evaluating their system on a set of 500 events with the highest association (similar to the event score by which we selected a set of best events). Their 37-class approach ultimately obtained a precision, recall and F1-score of 0.85, 0.55 and 0.67. Our system offered a comparable performance: a precision, recall and F1-score of 0.67, 0.67 and 0.65.

A comparable approach to predicting civil unrest is the EMBERS system by [16]. They evaluated their system over a period of a month, resulting in a precision and recall of respectively 0.69 and 0.82. In comparison, we obtained a higher precision while our estimated recall is lower. It is interesting to note how they received this recall score. They obtained a gold standard set of *social action events* by an independent organization that had human analysts survey newspapers and other media for mentions of civil unrest; arguably a reliable way of calculating recall in the real world. Our approach is only able to recall events that were present in the set of [8]. We have not explored ways to evaluate to what extent [8] detected all *social action events* that actually happened. We should consider the possibility that we might still miss *social action events* that were never detected as events in the first place, lowering our estimated recall.

Using the ranking of the Bayesian probabilities helped to increase the recall of *social action events* by 232%. We did not use the actual probabilities to influence the classification process, but used only the ranking of these probabilities. A potential direction for future research is to use the per-class probabilities generated by the Naive Bayes classifier in a more sophisticated manner. For example, it is possible to learn a certain probability threshold for *social action* and classify events that exceed this threshold as *social action*, regardless of the probability of other categories. The actual implementation of such a method requires a search for the best threshold setting. The main advantage of this approach is that this allows us to specify a specific precision-recall trade-off that is the most suitable for predicting social action events.

Our study has shown that the detection of *social action events* from Twitter based on open-domain event extraction and a subsequent event categorization

procedure is feasible. Due to the broad scope on open-domain events as starting point, we expect that this approach could be refined and improved when the focus is more on social action, e.g. by using lexicons of words associated with social action. Studying the extent of the potential added value of domain-specific knowledge is open for future work.

References

1. Bizer, C., Lehmann, J., Kobilarov, G., Auer, S., Becker, C., Cyganiak, R., Hellmann, S.: Dbpedia-a crystallization point for the web of data. *Web Semant. Sci. Serv. Agents World Wide Web* **7**(3), 154–165 (2009)
2. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent Dirichlet allocation. *J. Mach. Learn. Res.* **3**, 993–1022 (2003)
3. Compton, R., Lee, C.-K., Lu, T.-C., de Silva, L., Macy, M.: Detecting future social unrest in unprocessed Twitter data: emerging phenomena and big data. In: 2013 IEEE International Conference on Intelligence and Security Informatics (ISI), pp. 56–60. IEEE (2013)
4. De Smedt, T., Daelemans, W.: Pattern for Python. *J. Mach. Learn. Res.* **13**(1), 2063–2067 (2012)
5. Hayes, A.F., Krippendorff, K.: Answering the call for a standard reliability measure for coding data. *Commun. Methods Meas.* **1**(1), 77–89 (2007)
6. Korkmaz, G., Cadena, J., Kuhlman, C.J., Marathe, A., Vullikanti, A., Ramakrishnan, N.: Multi-source models for civil unrest forecasting. *Soc. Netw. Anal. Min.* **6**(1), 1–25 (2016)
7. Krippendorff, K.: *Content Analysis: An Introduction to Its Methodology*. Sage, Thousand Oaks (2004)
8. Kunneman, F., van den Bosch, A.: Automatically identifying periodic social events from Twitter. In: *Proceedings of the RANLP 2015*, pp. 320–328 (2015)
9. Kunneman, F., van den Bosch, A.: Open-domain extraction of future events from Twitter. *Nat. Lang. Eng.* **22**, 655–686 (2016)
10. Meij, E., Weerkamp, W., de Rijke, M.: Adding semantics to microblog posts. In: *Proceedings of the Fifth ACM International Conference on Web Search and Data Mining*, pp. 563–572. ACM (2012)
11. Mihalcea, R., Csomai, A.: Wikify!: linking documents to encyclopedic knowledge. In: *Proceedings of the Sixteenth ACM Conference on Information and Knowledge Management*, pp. 233–242. ACM (2007)
12. Muthiah, S., Huang, B., Arredondo, J., Mares, D., Getoor, L., Katz, G., Ramakrishnan, N.: Planned protest modeling in news and social media. In: *AAAI*, pp. 3920–3927 (2015)
13. NOS: Cohen: fouten politie, burgemeester. Nederlandse Omroep Stichting, 7 March 2013. <http://nos.nl/>
14. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: machine learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011)
15. Ramage, D., Dumais, S.T., Liebling, D.J.: Characterizing microblogs with topic models. *ICWSM* **10**, 1 (2010)

16. Ramakrishnan, N., Butler, P., Muthiah, S., Self, N., Khandpur, R., Saraf, P., Wang, W., Cadena, J., Vullikanti, A., Korkmaz, G., et al.: ‘Beating the news’ with embers: forecasting civil unrest using open source indicators. In: Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 1799–1808. ACM (2014)
17. Ritter, A., Etzioni, O., Clark, S., et al.: Open domain event extraction from Twitter. In: Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 1104–1112. ACM (2012)
18. van Heerden, D.: Facebook birthday invite leads to mayhem in Dutch town, authorities say. CNN, 24 September 2012. <http://edition.cnn.com/>
19. Volkskrant: Enkele duizenden bij protestmars bezuinigen. de Volkskrant, 21 September 2013. <http://volkskrant.nl/>
20. Zhao, W.X., Jiang, J., Weng, J., He, J., Lim, E.-P., Yan, H., Li, X.: Comparing Twitter and traditional media using topic models. In: Clough, P., Foley, C., Gurrin, C., Jones, G.J.F., Kraaij, W., Lee, H., Mudoch, V. (eds.) ECIR 2011. LNCS, vol. 6611, pp. 338–349. Springer, Heidelberg (2011). doi:[10.1007/978-3-642-20161-5_34](https://doi.org/10.1007/978-3-642-20161-5_34)