

Learning Topical Social Sensors

Authors

Affiliations

Abstract

Social media sources such as Twitter represent a massively distributed social sensor of a [kaleidoscope of topics](#) ranging from social and political events to entertainment and sports news. [We note, however, that due to the overwhelming volume of content, it can be difficult to spot novel and significant topics within a broad theme in a timely fashion – such as #obamacare or a recent #twochild policy in China.](#) This paper propose a scalable and practical method to automatically construct social sensors for generic topics. we train a supervised learner to identify topical content [from millions of features capturing content, user and social interactions on Twitter](#). On a corpus of approximately 1 billion English Tweets collected from the Twitter streaming API during 2013 and 2014 and learning for 10 diverse [themes](#) ranging from social issues to celebrity deaths to the “Iran nuclear deal”, we empirically show that our learned social sensor automatically generalizes to unseen future content with high ranking and precision scores. Furthermore, we provide an extensive analysis of features and feature types across different topics that reveals, for example, that (1) largely independent of topic, simple terms are the most informative feature followed by location features and that (2) the number of unique hashtags and tweets by a user correlates more with their informativeness than their follower or friend count. In summary, this work provides a novel, effective, and efficient way to learn topical social sensors from a seed set of hashtags requiring minimal user curation effort and offering strong generalization to future topical content.

1 Introduction

Social media sites such as Twitter present a double-edged sword for users. On one hand these sources contain a vast amount of novel and topical content that challenge traditional news media sources in terms of their timeliness and diversity. Yet on the other hand they also contain a vast amount of spam and otherwise low-value content for most users’ information needs where filtering out irrelevant content is extremely time-consuming. Hence, while it is widely acknowledged that social media sources can be used as topical content sensors (indeed, an entire European Union project was focused on related “Social Sensor” research¹), automati-

cally learning high-precision sensors (i.e., ranking and retrieval methods) for arbitrary topics that generalize to future unseen content remains an open question in the literature and comprises the key problem we seek to address in this paper.

Perhaps the critical bottleneck for learning targeted topical social sensors is to achieve sufficient supervised content labeling. With data requirements often in the thousands of labels to ensure effective learning and generalization over a large candidate feature space (as found in social media), manual labeling is simply too time-consuming for many users and crowdsourced labels are both costly and prone to misinterpretation of users’ information needs. Fortuitously, hashtags have emerged in recent years as a pervasive topical proxy on social media sites — hashtags originated on IRC chat, were adopted later (and perhaps most famously) on Twitter, and now appear on other microblogs (e.g., Sina and Tencent Weibo) and even Facebook. Hence as a simple enabling insight that serves as a catalyst for effective topical social sensor learning, we leverage a (small) set of user-curated topical hashtags to efficiently provide a large number of supervised topic labels for social media content.

With the data labeling bottleneck resolved, we proceed to train supervised classification and ranking methods to learn topical content from a large feature space of source users and their locations, terms, hashtags, and mentions. On a corpus of approximately 1 billion English Tweets collected from the Twitter streaming API during 2013 and 2014 and covering 10 diverse topics ranging from social issues to celebrity deaths to the “Iran nuclear deal”, we empirically show that two simple and efficiently trainable methods — logistic regression and naive Bayes — generalize well to unseen future topical content (including content with no hashtags) in terms of their mean average precision (MAP) and Precision@ n for a range of n . Furthermore, we show that terms and locations are among the most useful features — surprisingly more so than hashtags, even though hashtags were used to label the data. And perhaps even more surprisingly, the number of unique hashtags and tweets by a user correlates more with their informativeness than their follower or friend count.

In summary, this work fills a major gap in the literature of topical social sensors and how to effectively and efficiently learn them given minimal supervision from a user. Our results suggest that these sensors generalize well to unseen future topical content and provide a novel paradigm for the

extraction of high-value content from social media.

2 Dataset Statistics

We crawled Twitter data using Twitter Streaming API for two years spanning 2013 and 2014 years. The total number of tweets collected is 829, 026, 458. In the context of Twitter, we consider a list of 5 features for each tweet. Each tweet has a *From*, the person who tweeted it, and a *Time* which is the date information of the tweet. It can also contain

- *Hashtag(s)*, keywords specified using # sign
- *Mention(s)*, another Twitter username being mentioned using @ sign
- *Term(s)*, uni-grams which we extract from the 140 characters of the tweet.

We provide more detailed statistics about each feature in Table 1 and Table 2. Here, we provide the unique number of the feature in our dataset, in addition to maximum, average, and median values of each feature across the tweets, user, and hashtag dimensions. For example, a hashtag has been used in average by 10.08 users or users have used 2 hashtags on average.

For *Location* feature, Fig 1 shows the distribution of tweets across different U.S. and international locations for 3 chosen topics which we found to be more interesting due to their geographical distribution over various locations. For example, we can see that Middle east and Malaysia stand out for the topic of Human Caused Disaster. Malaysia is a hot place for this topic due to MH370 incident with lots of usage of #whereisthefuckingplane.

Used in #Tweets				
Feature	Max	Avg	Median	Max entity
From	10,196	8.67	2	running_status
Hashtag	1,653,159	13.91	1	#retweet
Mention				
Location				
Term	241,896,559	492.37	1	rt
Used by #Users				
Hashtag	592,363	10.08	1	#retweet
Mention	26,293	5.44	1	dimensionist
Location	739,120	641.5	2	london
Term	1,799,385	-	1	rt
Using #Hashtags				
From	18,167	2	0	daily_astrodata

Table 1: Feature Statistics of 829, 026, 458 tweets in our Twitter dataset

	From	Hashtag	Mention	Location	Term
#Unique	95,547,198	11,183,410	411,341,569	58,601	20,234,728

Table 2: Number of unique values for each feature of 829, 026, 458 tweets in our Twitter dataset

3 Experimental Methodology

With the 5 set of features defined as *From*, *Mention*, *Location*, *Term*, *Hashtag*, we proceed to define the methodology for retrieving ranked list of tweets for a given topic. The list of topics were defined to be a set of 10 various topics covering very specific e.g., *IranDeal*, and very broad e.g., *SocialIssues* topics: *Tennis*, *Space*, *Soccer*, *IranDeal*, *HumanDisaster* (HumanCausedDisaster), *CelebrityDeath*, *SocialIssues*, *NaturalDisaster*, *Epidemics*, and *LGBT*.

Our goal is to retrieve a ranked list of tweets T_i by employing machine learning methods using defined features. Fig ?? showed unique number of values for each feature. These values sum up to a total number of 538,365,507 features. Also, as noted earlier, we are dealing with 829,026,458 number of tweets. This shows the need for providing techniques to:

1. Annotate the tweets as topical and non-topical
2. Select a set of features for learning the model

In regards to annotating the tweets, first, a set of topical hashtags are manually curated for each topic. This set is annotated manually with 2 annotator individually and inner-annotator agreement was achieved by reviewing these sets by 2 more individuals. Each *Hashtag* has a birthday which is defined as the first time it has been used in our dataset. The criteria of choosing hashtags included

- To be related to the topic
- To be preferably born during our time span i.e. not being used before e.g., #ebola
- The entire set of hashtag birth dates to cover the two years of our data

After choosing these sets, we tag a tweet as topical if it contains at least one topical hashtag.

In order to conduct our experiments, tweets are temporally divided over 2 years to provide train and test sets. Since our tweet labeling is through topical hashtags, this division is done in a way to retain enough hashtags for train, validation, and test timespan. To this purpose, hashtags are divided based on their birth date with 50 percent of hashtags being born at train timespan, 10 percent born at validation timespan, and the last 40 percent born at test timespan. Table 3 provides samples of hashtags, number of train hashtags, test hashtags, and topical tweets for each topic. We can see that some topics such as *HumanDisaster* and *Soccer* are more general topics and have higher number of topical tweets while some other ones such as *IranDeal* is more specific, thus having less number of topical tweets.

Regarding feature selection, it is clear that it is not possible to learn a model with total number of 538,365,507 features. Even if it was possible, we would have a problem for providing enough training samples, and our feature vectors would be extremely sparse considering 140 characters limitation of Twitter. Therefore, we performed a primary feature selection based on frequency of each feature. The feature selection process included:

- Cleaning *Term* feature to remove stop-words

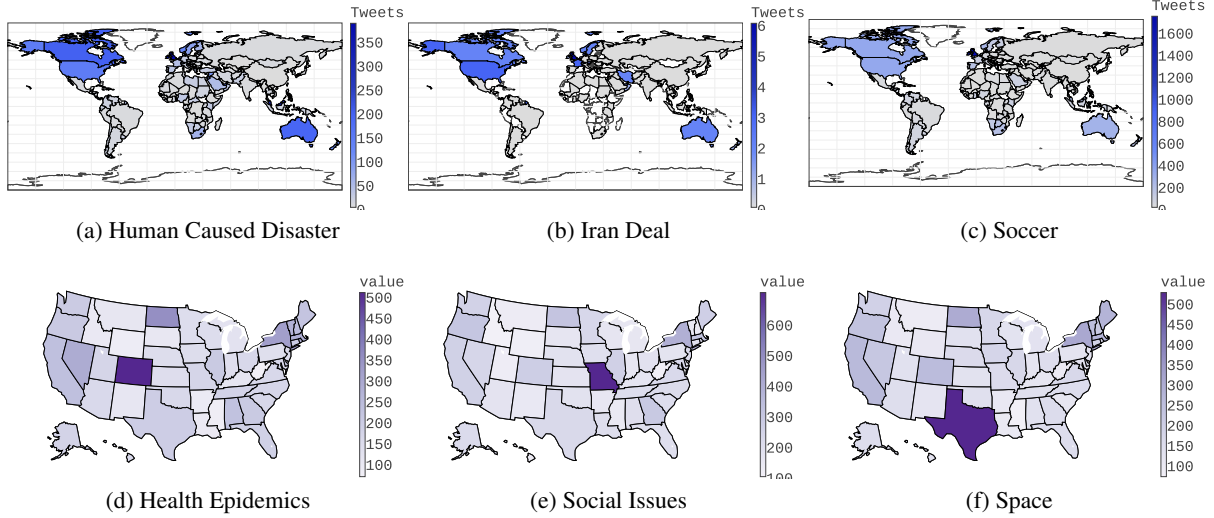


Figure 1: Distribution of tweets across International locations (top row) and U.S. locations (bottom row)

- Choosing a cut-off threshold of 159 for *From*, *Mention*, *Term* features
- Choosing a cut-off threshold of 50 for *Location* and *Hashtag* features (These features had much lower number of unique values)

This results in roughly 1 million features, denoted as social features set SF where

$$\begin{aligned}
 SF_m \in & \{Hashtag_{m1}, From_{m2}, Mention_{m3}, Term_{m4}, Location_{m5}\} \\
 m = & \{1, \dots, 1166582\}, m1 = \{1, \dots, 184702\}, \\
 m2 = & \{1, \dots, 361789\}, m3 = \{1, \dots, 244478\}, \\
 m4 = & \{1, \dots, 317846\}, m5 = \{1, \dots, 57767\} \quad (1)
 \end{aligned}$$

Classification Algorithms

Now that we defined the primary steps for preparing the features and dataset, we can use them to build a ranking approach for topical tweet selection. Our method is based on classification/ranking approaches defined in the literature to weight features. These weights are further used to rank tweets for each topic. Here, we use the following classification approaches:

1. Logistic Regression
2. Naive Bayes
3. Rocchio (centroid)
4. RankSVM

To this purpose, we define the problem as assigning a weight W_i to tweet X as a measure of similarity to the given topics t_i . W_i is the sum of weights of features in the x_i :

$$W_i = \sum_k w_k \times f_k \quad (2)$$

where w_k is the weight of feature f_k and $f_k \in \{true, false\}$ represents whether each of the features in $SF_m, m = \{1, \dots, 1166582\}$ is present in tweet X or not. The weights w_k are learned by applying one of the classification algorithms.

For example, in case of Naive Bayes, the problem is defined as probability of tweet being in topic t_i given feature vector F_x for

$$P(t_i|X) = P(t_i|F_x) = \frac{P(t_i)P(F_x|t_i)}{P(F_x)} \quad (3)$$

In order to learn the models, we take the following steps for each topic:

1. All the positive tweets for the given topic, in addition to sub-sampled set of negative tweets are collected
2. A set of top N features are selected based on the Mutual Information values of features for the given topic. N is selected during hyper-parameter tuning phase from the set of $10E1, 10E2, 10E3, 10E4, 1166582$ values.
3. Train, validation, and train sets are further built based on the division process explained in 3 and selected set of top N features.
4. The model's hyper-parameter in addition to N parameter are tuned on the validation set
5. The model learns the weight vector based on the tuned hyper-parameters on the full train set and validation set

The Liblinear (?) package is used for implementing *LR* and *RankSVM*. The reason for deciding to tune the models on top N features based on Mutual Information, comes from our primary feature analysis on the dataset which showed the ability of Mutual Information measure to pick more correlated features for each topic. This is discussed in more details in section 4. The model hyper-parameters are tuned for *LR* and *NB*. The *Rocchio* method is parameter free and

	Tennis	Space	Soccer	IranDeal	HumanDisaster	CelebrityDeath	SocialIssues	NaturalDisaster	Epidemics	LGBT
#TrainHashtags	58	98	126	12	49	28	31	31	52	29
#TestHashtags	36	63	81	5	29	16	19	19	33	17
#TopicalTweets	55,053	239,719	860,389	8,762	408,304	163,890	230,058	230,058	210,217	282,527
Sample Hashtags	#usopenchampion	#asteroids	#worldcup	#irandeal	#gazaundersattack	#robinwilliams	#policebrutality	#policebrutality	#ebola	#loveislove
	#novakdjokovic	#astronauts	#lovesoccer	#iranfreedom	#childrenofsyria	#ripmandela	#michaelbrown	#michaelbrown	#virus	#gaypride
	#wimbledon	#satellite	#fifa	#irantalk	#iraqwar	#rijoanivers	#justice4all	#justice4all	#vaccine	#uniteblue
	#womenstennis	#spacecraft	#realmadrid	#rouhani	#bombthreat	#mandela	#freetheweed	#freetheweed	#chickenpox	#homo
	#tennisnews	#telescope	#beckham	#nuclearpower	#isis	#paulwalker	#newnjgunlaw	#newnjgunlaw	#theplague	#gaymarriage

Table 3: Test/Train Hashtag samples and statistics

the LibLinear (?) implementation of *RankSVM* does not provide manual tuning of the model’s hyper-parameter.

Analysis

After experimenting each mentioned model on our dataset, we provide the following metrics:

- MAP: Mean average precision for a set of topics is the mean of the average precision scores for each topic.
- P@K: Precision at K for $K \in \{10, 100, 1000\}$, the number of relevant results on the first K search results page

The model’s hyper-parameters are tuned based on MAP scores, having MAP as our most important metrics. Table 4 provides these metrics for each topic. Logistic Regression method is the method that performs best on average. Generally, Naive Bayes performed comparable/better to Logistic Regression having second best average value of MAP. We also provide the top 5 tweets returned by Logistic Regression for each topic as anecdotal results in Table 5. In this table, the signs in the beginning of the tweet represent the following:

- ✗ represents the tweets that are method has incorrectly ranked as highly topical
- ✓ represents the tweets correctly ranked as highly topical
- ★ represents the tweets that don’t have any topical hashtags and therefore are not labeled as correctly ranked topical. However, looking at the tweets, we can see that they are in fact related to the topic

The fact that there are cases of tweets not being correctly labeled as topical, provides evidence that our method of labeling tweets has limitations and our MAP and P@K values are actually suffering from this problem. However, this shows the power of Logistic Regression method in generalizing from a small set of hashtags.

4 Feature Analysis

. In this section, we analyze the informativeness of each feature for learning topical tweets by looking at different characteristics for each feature in our dataset. For example, one characteristic of hashtags could be the number of the tweets that contain those hashtags. Does this have an effect on importance of the hashtag when it comes to learning topical tweets or not. In this sense, this section would bring insights to the following questions:

- What are the best features for learning social sensors, do they differ by topic? (Why?)

- For each feature type, do any attributes correlate with importance?

A famous method for measuring informativeness is Mutual Information which is a measure of amount of information one random variable contains about another random variable. In order to calculate amount of information that a feature $f_k \in \{from, hashtag, mention, term, location\}$ provides w.r.t $t_i \in \{NaturalDisaster, Epidemics, \dots\}$, mutual information is defined as:

$$I(t_i, f_k) = \sum_{t_i \in \{true, false\}} \sum_{f_k \in \{true, false\}} p(f_k, t_i) \log \left(\frac{p(f_k, t_i)}{p(f_k)p(t_i)} \right) \quad (4)$$

Higher values for this metric indicates more informative features for the specified topic.

In order to answer the first question on what are the best features for learning social sensors, we provide mean of Mutual Information values for each feature across different topics in Table 3. The last column in this table shows average of mean Mutual Information for the feature. The following observations are from the analysis of Table 3:

- *Term* feature is the most prevalent feature and in general, the more features you have, the better the chance that one is useful.
- *Location* and *Hashtag* feature provide second and third most informative features respectively.
- A few topics such as *IranDeal* and tennis are less sensitive to selection of a specific features.
- Location feature provides more information regarding *HumanDisaster*, *LBGT*, and *Soccer* topics.
- Sorting features based on their average mean values across different topics results in the following order: *Term*, *Location*, *Hashtag*, *Mention*, *From*

In general, this presents evidence on the need for learning the weights of features for each topic, because there is no specific selection of features that would separate various topics from each other.

Also, in order to show the power of Mutual Information criteria, we present the top 5 features for each topic in table 6. It can be observed how different locations, hashtags, or terms showed as the top features based on mutual information are actually in relation with the specific topic.

		Tennis	Space	Soccer	IranDeal	HumanDisaster	CelebrityDeath	SocialIssues	NaturalDisaster	Epidemics	LGBT	Mean
LR	MAP	0.918	0.870	0.827	0.811	0.761	0.719	0.498	0.338	0.329	0.165	0.623±0.19
NB	MAP	0.908	0.897	0.731	0.824	0.785	0.748	0.623	0.267	0.178	0.092	0.605±0.22
Rocchio	MAP	0.690	0.221	0.899	0.584	0.481	0.253	0.393	0.210	0.255	0.089	0.407±0.18
RankSVM	MAP	0.702	0.840	0.674	0.586	0.603	0.469	0.370	0.248	0.136	0.082	0.471±0.18
LR	P@10	1.000	0.000	0.200	0.700	0.600	0.000	0.100	0.200	0.300	0.500	0.360±0.24
NB	P@10	1.000	0.900	0.700	0.600	0.600	0.700	1.000	0.100	0.400	0.100	0.610±0.23
Rocchio	P@10	0.800	0.000	1.000	0.900	0.000	0.000	0.000	0.500	0.500	0.100	0.380±0.29
RankSVM	P@10	1.000	0.800	0.600	0.800	0.400	0.300	0.000	0.100	0.000	0.200	0.420±0.26
LR	P@100	0.950	0.580	0.650	0.870	0.620	0.490	0.640	0.690	0.790	0.210	0.649±0.15
NB	P@100	0.980	0.850	0.600	0.880	0.750	0.860	0.730	0.230	0.090	0.190	0.616±0.23
Rocchio	P@100	0.980	0.000	1.000	0.690	0.170	0.000	0.280	0.170	0.680	0.120	0.409±0.28
RankSVM	P@100	0.730	0.720	0.310	0.700	0.880	0.440	0.480	0.340	0.020	0.100	0.472±0.20
LR	P@1000	0.963	0.954	0.816	0.218	0.899	0.833	0.215	0.192	0.343	0.071	0.550±0.26
NB	P@1000	0.954	0.954	0.716	0.218	0.904	0.881	0.215	0.195	0.141	0.060	0.524±0.28
Rocchio	P@1000	0.604	0.000	0.925	0.218	0.359	0.000	0.215	0.167	0.144	0.065	0.270±0.21
RankSVM	P@1000	0.799	0.922	0.764	0.218	0.525	0.547	0.215	0.173	0.154	0.064	0.438±0.22

Table 4: Different learning methods results on topics with hyper-parameter tuning based on MAP

Tennis	Space
✓ rt @espnentemis: shock city. darcis drops rafa in straight sets. first time nadal loses in first rd of a. major in career. #espnwimbledon #w	✗ rt @jaredleto: rt @30secondstomars: icymi: mars performing a cover of @rihanna's #stay on australia's @triplemelb - video - http://t.co/uq
✓ @ESPN/Tennis: Shock city. Darcis drops Rafa in straight sets. First time Nadal loses in first rd of a. Major in career.	✗witing mars @30secondstomars @jaredleto @shannonleto @tomfromearth vibest group http://t.co/dloovrjiaf
✓ @ESPN/Tennis: Djokovic ousts the last American man standing @Wimbledon, beating Reynolds 7-6 6-3 6-1 #ESPNWimbledon	✗rt @jaredleto.com: show everyone how much you are proud of @30secondstomars. #mtv/hottest 30 seconds to mars http://t.co/byxni467
✓ Nadal's a legend. After 3 years; Definitely He's gonna be the best of all the time. Unbelievable performance. @RafaelNadal #USOpenFinal	✗rt @30secondstomars: missed the big news? mars touring with @linkinpark + special guests @afi this summer!. http://t.co/3e5m0pwrld
✓ @calvy70 @ESPN/Tennis @Wimbledon I see, thanks for the info and enjoy #Wimbledon2014	✗rt @30secondstomars: to the right, to the left, we will fighto the death.go #intothewildonyrt with mars, starting weekly, nov 30 . hit
Soccer	IranDeal
✗rt @tomm_dogg: #thingsdoforeartheners spend all my money.	✓ rt @iran_policy: @vidalquadrax: @isjcommittee has investigated 10 major subjects of irans controversial #nuclear program #irantalksvienna
✗rt @mancionlineco nice performance	✓ rt @iran_policy: @vidalquadrax: @isjcommittee has investigated 10 major subjects of irans controversial #nuclear program #irantalksvienna
✗rt @indykalla: podolski: "let's see what happens in the winter. the fact is that i'm not happy with it, that's clear." @arsenal	✗rt @negarmortazavi: thank you @hassanrouhani for retweeting. let's hope for a day when no iranian fears returning to their homeland. http://
✗rt @indykalla: Wenger: "i don't believe match-fixing is a problem in england." #afc	✗rt @iran_policy: iran: details of savage attack on political prisoners in evin prison http://t.co/sdzaukqdiv #iran #humanrights
✗rt @indykalla you never got back to me about tennis this week	✓ rt @iran_policy: chairman ros-lehtinen speaking on us commitment 2 protect camp liberty residents. #iranhrviolations http://t.co/lg6dhlznu
HumanDisaster	CelebrityDeath
✓ rt @baselsyrian: there've been peaceful people in #homs not terrorists! #assad.enemy of #humanity destroyed it. #eyehomhs #withsyria http:	✗rt @sawubona_chris: today is my birthday & also the day my hero @nelsonmandela has died. lets never forget what he taught us. forgiveness i
✓ what a helpless father, he can do nothing under #assad's siege! #speakup4syrianchildren http://t.co/vglc3byebw#syria #syriawarcrimes #un	✗rt @nelsonmandela: death is something inevitable when a man has done what he considers to be his duty to his people&#amp;this country,he can res
✗exclusive: us formally requested #un investigation; russia pressured #assad to no avail;chain of evidence proof hard http://t.co/5602rvidfw	✗rt @nelsonmandela: la muerte es algo inevitable cuando un hombre ha hecho lo que considera que es su deber para con su gente y su pas,pued
✗#save_aleppo from #assadwarcrimes#save_aleppo from #civilians -targeted shelling of #assad regime#syria #aleppo http://t.co/k3dixh0pxl	✗rt #jacques #kallis: a phenomenal cricketing giant of all time - #cricket #history #southafrica http://t.co/ms5pmwoug9
✓ rt @camille_rhites: why does the #un allow this to continue? rt @tintin1957 help raise awareness of the suffering in #syriawarcrimes http://t	✗rt @sudehs1304 south africa has the most beautiful babies...so diverse,so unique...so god!! lol #durban #southafrica
SocialIssues	NaturalDisaster
✗the us doesn't actually borrow is the thing. i believe in a creationist theory of the us dollar @usanationdebt @nationaldebt	✗us execution in #oklahoma : not cruel and unusual? maybe just barbaric, inhumane and reminiscent of the dark ages!
✗rt @2anow: according to @njsenatpex women's rights do not include this poor nj mother's right to defend herself http://t.co/czbslqkht6 #	✗rt #haiti #politics - the haiti-dominican crisis - i agree with how martelly is handling the situation: i totally... http://t.co/ro4pwwsxs
✗rt @2anow: confiscation ? how many carry permits are in the senate and assembly? give us ours or turn them in. @senatororettaw @lougreenw	✗rt @soilhaiti: a new reforestation effort in #haiti. local compost, anyone? http://t.co/spadrtqbjk @richardbrannon @clintonfdn @virginunite
✗rt @2anow: vote with your wallet against #guncontrolforest city enterprises does not support the #2a http://t.co/tpk0k3herm#nj2as #tco	✗mex cousins jamais ns hantent les nuits de duvalier #haiti #duvalier
✗rt @2anow @momsdemand @jstines3 they don't have a plan for that, which is why they should never be allowed to take our guns	✓ tony burgener of @swissolidarity says you can't compare the disaster response in #haiti with the response to #hayan in #philippines @heid
Epidemics	LGBT
✓ rt @who: fourteen of the susp. & conf. ebola cases in #conakry, #guinea, are health care workers, of which 11 died #askebola	✗rt @jackmoldcuts: @lunaticrex @fingersmalloy @toddkincannon @theanonliberal anthony kennedy just wrote opinion granting...
✗rt @who who can afford also been cover in government health insurance [with universal health coverage]	✗rt @toddkincannon your personal account, your interest, separate from your business.
✓ #ebolabreak this health crisis...unparalleled in modern times, @who dir. alyward - requires \$1 billion to stem http://t.co/rjzqhydh3d	✗why would you report someone as spam if he is not spam? @ilgyirbrea @toddkincannon
✗rt @medsin: @who are conducting a survey on the social determinants of health in medical teaching. fill the survey in at https://t.co/aj59x	✗rt @th.arch3r: @toddkincannon thanks for your lt having the female realbrother. between them is 600 lbs. 104 iq points. and a lot of hate.
✗augmentation vertigineuse de 57,4% en 1 an des actes islamophobes en france, dit le collectif contre l'islamophobie http://t.co/2qihocge15	✗rt @toddkincannon who us dick trickle.

Table 5: Top Tweets for each topic based on MAP tuned results

In order to answer the second question on whether any attributes correlate with importance for each feature, we provide two set of analysis. The first one, provides Mutual Information values of each feature across feature's attribute values shown by violin plots in figure 4. The attributes for each feature are:

- From: favorite count (the number of tweets the user has favorited), followers count (the number of users who follow the user), friends count (the number of users followed by the user), hashtag count (number of hashtags used by the user), tweet count (the number of tweets from the user)
- Hashtag: tweet count, user count (the number of users using the hashtag)
- Location: user count
- Mention: tweet count
- Term: tweet count

As we can see in the violin plots, the general pattern is that the more number of tweets, users, or hashtags count a feature has, the higher the chance of becoming topical will be. This

pattern exists on other attributes of *From* feature, although a bit less clear than the tweets, users, or hashtags counts attributes. In addition, we further analyzed the density plots of favorite count, follower count, friends count, hashtag count attributes of *From* feature shown in Fig 5. These plots represent a bi-modality in the distribution. Further analysis of data showed that the top mode belongs to users who have at least one topical tweet while bottom mode are users with no topical tweet.

5 Related Works

This section documents existing research on the use of social media as a sensor for topic detection on social media. Herein, we focus on related research on both events and topics detection within social media. With the consideration that events are special type of topics and can be classified as such. To see how different works address topic detection on social media, we focus on three extensively researched types of topic detection: trending topic detection, specific event detection, and tweet recommendation.

The first overarching group of works reviewed herein fo-

Topics/Top10	NaturalDisaster	Epidemics	IranDeal	SocialIssues	LGBT	HumanDisaster	CelebrityDeath	Space	Tennis	Soccer
From	earthquake_wo	changedecopine	mazandara	nsingerdebtpaid	eph4.15	ydumozyf	nmandelaquotes	daily_astrodata	tracktennisnews	losangelessrh
From	earthalerts	drdaveanddee	hhadi119	debtadvisoruk	mgdauber	syriatweeten	boiknox	freesolarleads	tennis_result	shootale
From	seelites	joinmentormetwk	140iran	debt_protect	stevendickinson	tintin1957	jacanews	houston_jobs	i_roger_federer	sport_agent
From	globalfloodnews	followebola	setarehgan	negativeequityf	lileensvf1	sirajsol	ewnreporter	star_wars_gifts	tennislessnow	books_you_want
From	gcmcdrought	localnursejobs	akhgarshabaneh	dolphin Js	truckerbooman	rt3syria	paulretweet	lenautilus	kamranisbest	makeupbella
Hashtag	earthquake	health	iran	ferguson	tcot	syria	rip	science	wimbledon	lfc
Hashtag	haiyan	uniteblue	irantalks	mikebrown	p2	gaza	riprobinwilliams	starwars	usopen	worldcup
Hashtag	storm	ebola	rouhani	ericgarner	pjnet	isis	ripcorymonteith	houston	tennis	arsenal
Hashtag	tornado	healthcare	iranian	blacklivesmatter	uniteblue	israel	mandela	sun	nadal	worldcup2014
Hashtag	prayforthephilippines	depression	no2rouhani	fergusondecision	teaparty	mh370	nelsonmandela	sxsw	wimbledon2014	halamadrid
Location	philippines	usa	tehran	st.louis	usa	malaysia	southafrica	germany	london	liverpool
Location	ca	ncusa	u.s.a	mo	bordentown	palestine	johannesburg	roodepoort	uk	manchester
Location	india	garlandtx	nederland	usa	newjersey	syria	capetown	houston	india	london
Location	newdelhi	oh-sandiego	iran	dc	sweethomealabama!	israel	pretoria	austin	pakistan	nigeria
Location	newzealand	washington	globalcitizen	washington	aurora	london	durban	tx	islamabad	india
Mention	oxfamgb	foxtramedia	4freedomiran	deray	jjauthor	ifalasteen	nelsonmandela	bizarro.chile	wimbledon	lfc
Mention	weatherchannel	obi.obadike	iran_policy	natedrug	2anow	revolutionsyria	realpaulwalker	nasa	usopen	arsenal
Mention	redcross	who	hassanrouhani	antoniofrench	govchristie	drbasselabuward	robinwilliams	j_ksen	andy_murray	realmadriden
Mention	twcbreaking	obadike1	un	bipartisanism	a5h0ka	mogaza	rememberrobin	jaredleto	serenawilliams	ussoccer
Mention	abc7	c25kfree	statedept	theanmessage	barackobama	palestinianism	tweetlikegiris	30secondstomars	esptennis	mcfc
Term	philippines	health	iran	police	obama	israel	robin	cnblue	murray	madrid
Term	donate	ebola	regime	protesters	gun	gaza	williams	movistar	tennis	goal
Term	typhoon	acrxx	nuclear	officer	rights	israeli	nelson	enero	federer	cup
Term	affected	medical	iranian	protest	america	killed	mandela	imperdible	djokovic	manchester
Term	relief	virus	resistance	cops	gop	children	cory	greet	nadal	match

Table 6: Top 5 features for each topic based on Mutual Information

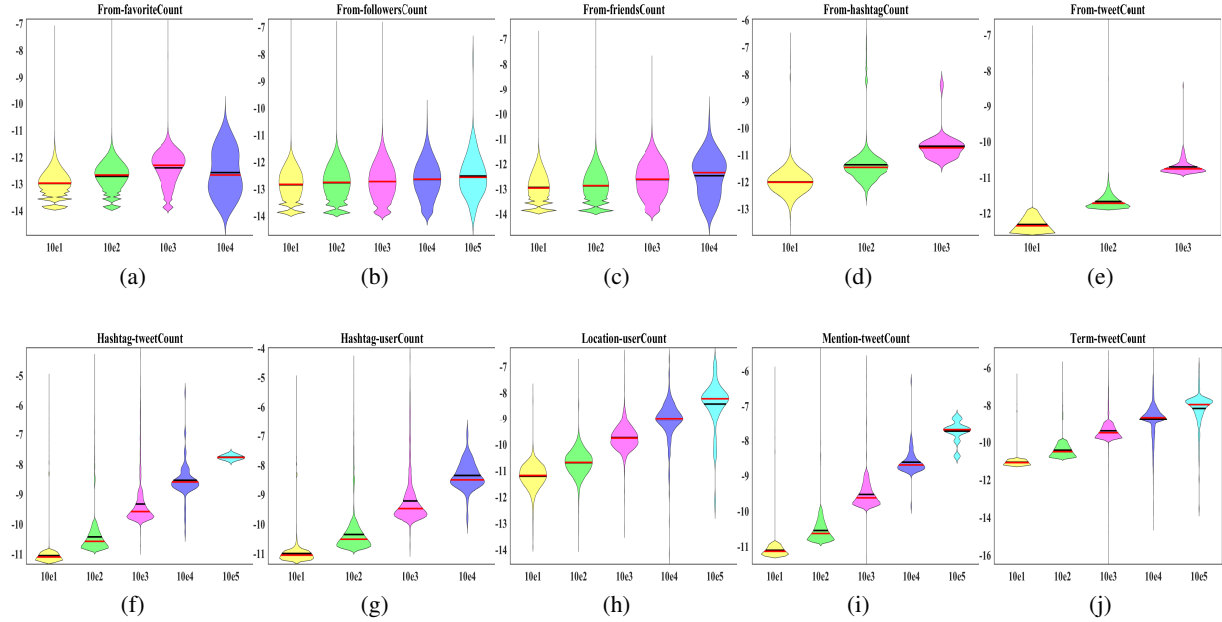


Figure 4: ViolinPlots for feature attributes counts vs. MI. Top row shows attributes {favoriteCount, followerCount, friendCount, hashtagCount, tweetCount} for *From* feature. Bottom row shows attributes tweetCount and/or userCount for *Hashtag*, *Location*, *Mention*, and *Term* features.

cus on trending topic detection methods. The majority of works detecting trending topics use bursts as the indicator of events, where a burst is defined as a sudden change in posting rates of some keywords, hashtags, etc. These can further be divided into multiple categories based on how they use bursts to extract the event. The first category, clustering-based methods, focuses on the hypothesis that trends are topical and topics are defined by the collection of relevant content; hence trends can be detected by

clustered content (Petrović, Osborne, and Lavrenko 2010; Ishikawa et al. 2012; Phuvipadawat and Murata 2010; Becker, Naaman, and Gravano 2011; O'Connor, Krieger, and Ahn 2010; Weng and Lee 2011). With more focus on machine learning methods, (Wei et al. 2015) proposed a graphical model to discover latent events clustered in the spatial, temporal and lexical dimensions, while (Yamamoto and Satoh 2015) focused on the task of multi-label classification of tweets into living aspects such as eating. The sec-

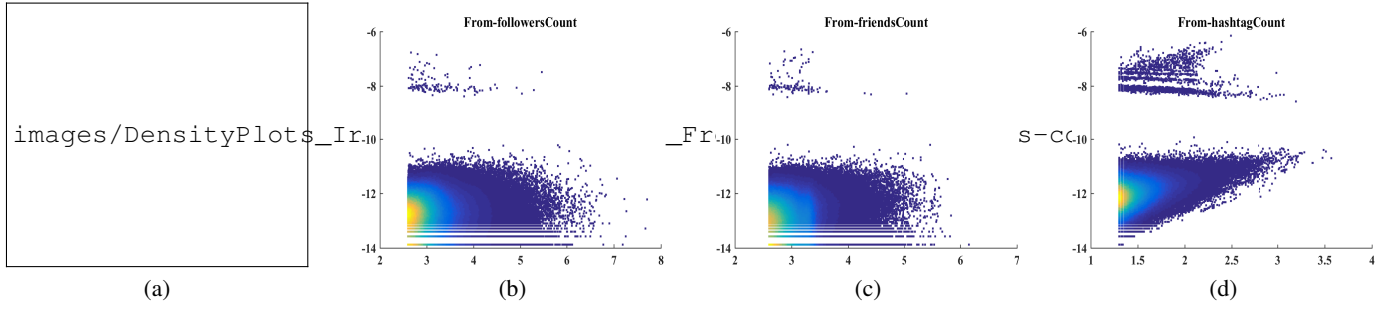


Figure 5: DensityPlots for feature attributes counts vs. MI. (a-d) show attributes {favoriteCount, followerCount, friendCount, hashtagCount} for *From* feature

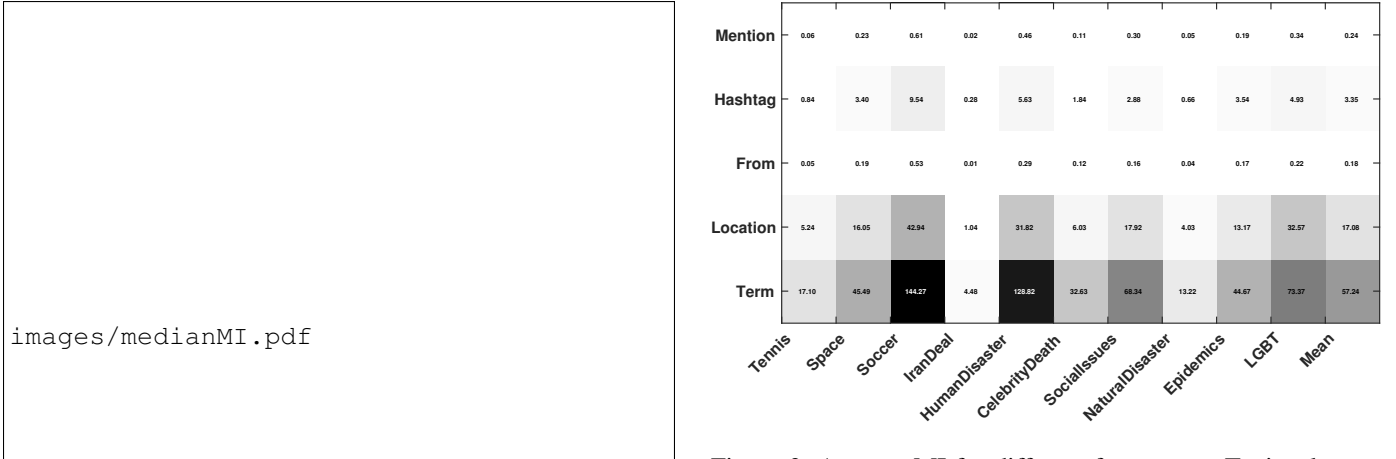


Figure 3: Average MI for different features vs. Topics, last two column show mean value and stderr across all topics

Figure 2: Median MI for different features vs. Topics, last two column show mean value and stderr across all topics

ond category, term-based methods focuses on the hypothesis that topics can be detected by focusing on temporal patterns of terms/keywords independent of the content of documents (Mathioudakis and Koudas 2010; Cui et al. 2012; Zhao et al. 2011; Nichols, Mahmud, and Drews 2012). The third category, query-based methods, focuses on the hypothesis that trending topics can be detected by measuring user defined criteria (Albakour, Macdonald, and Ounis 2013; Sakaki et al. 2012). The fourth category, network Structure-based method, focuses on the hypothesis that trending topics can be detected by studying the network structure of users (Budak, Agrawal, and El Abbadi 2011). The final category, hybrid method of (Diplaris et al. 2012) introduced concept of Dynamic Social Containers in this work to take advantage of aggregation of mining both the structure, content, and

multimedia data to index and provide personalized, context-aware search. In this work, the authors defined social sensor as analyzing the dynamic and massive amount of information provided by user with the purpose of extracting unbiased trending topics and events in addition to using social connections for recommendation.

With the purpose of comparison of methods, (Aiello et al. 2013) evaluated six trending topic detection methods on three Twitter datasets differing in time scale and topic churn rate. The authors conclude that natural language processing techniques perform well on focused topics. However, techniques mining temporal distribution of concepts are needed to handle more heterogeneous streams.

However, trending topics detection methods are not targeted. Our method differs from trending topic detection methods in that we are focusing on a set of topics that cannot necessarily be detected using bursts. Thus, trending topics detection methods are of limited relevance to the work presented hereinafter.

The second overarching group of works focuses on detection of a specific targeted topic, such as a disaster or epidemic. In a predictive study by (Kryvasheyev et al. 2014), the authors studied the network of users and focused on choosing the best groups of users in order to achieve lead-

times i.e. faster detection of disastrous event (following the concept of "friendship paradox" (Feld 1991)²). (Sakaki, Okazaki, and Matsuo 2013) used SVM classifier to detect earthquakes and employed a location estimation method such as Kalman Filtering for localizing it. The authors detected the occurrence of earthquakes through extracted statistical features e.g., the number and position of words in a tweet, keyword features and word context features from tweets.

Whereas the above works addressed exploiting the detection of crisis events, the following works focused on descriptive studies on disaster. The studies discuss the behavior of Twitter users during a crisis (Vieweg et al. 2010; Cheong and Cheong 2011; Starbird and Palen 2010) and do not address exploiting detection of crisis events. The studies investigated the use of social media during a crisis in order to identify information propagation properties, the social behavior of users (their retweeting behavior), information contributing to situational awareness, and the active players in communicating information. The behavioral information gleaned from these studies is exploited in this work to aid in the development of social sensors for detection of topics.

To detect health epidemics, researchers used content-based and/or structure-based methods. The content-based methods of (Culotta 2010) and (Aramaki, Maskawa, and Morita 2011) identified influenza-related tweets and correlated these tweets to United States Center for Disease Control (CDC) statistics on influenza, such as the infection and incubation rate. As for methodology, both works extracted bag-of-words as features, while the former employed single and multiple linear regression showing that multiple linear regression works better, while the latter employed SVM. Results indicated a high correlation between their estimation of influenza cases in early stages of an epidemic, and statistics from the CDC and Japan's Infection Disease Surveillance Center. The other approach to early detection of contagious outbreaks is to use structure-based methods, (García-Herranz et al. 2012) designed a sensor based on the friendship paradox concept for early detection of contagious outbreaks. In this regard, García-Herranz et al. provided a method for choosing sensor groups from friends of random sets of users to find more central individuals in order to enforce early detection. The central assumption made in this work is that a sensor group represents more central individuals, and individuals at the center of a network are more likely to become infected than randomly-chosen members of the population. As a result, (García-Herranz et al. 2012) argued that this selection process of sensor groups helps in the early detection of outbreaks.

On the other hand, hybrid method of (Sadilek, Kautz, and Silenzio 2012), exploited tweet content and the structural information of a user's network. The authors employed a semi supervised approach to learn a SVM classifier, using n-grams as features in order to detect ill individuals. Using co-location and friendship, the authors estimated the probability of physical interaction between healthy and sick peo-

²On average, most people have fewer friends than their friends have

ple. This enabled them to study the effect of these two factors of social activity (co-location for contact network and friendship for social ties) on public health.

The limitations of these studies centers on the fact that the proposed methods are only valid for detecting a single topic. These methods used a primitive methods for curating the data e.g., querying keyword earthquake. In addition, there is no discussion within these works on how these methods can be generalized for other topics.

Another set of studies have moved towards creating more generalizable methods. Using a dataset of 55,000 news articles and 121,000 tweets, (Krestel et al. 2015) compared four different methods of language model, topic model, logistic regression, and boosting, to evaluate recommended tweets for a given news article.. (Yan, Lapata, and Li 2012; Chen et al. 2012) also focused on tweet recommendation. Their methods considered the users twitter profile, including tweet and retweet history, and social relations as features. Coupled with tweet popularity, the methods are able to generate tweet recommendations. With the purpose of photo recommendation on social media websites, (Chiarandini et al. 2013) analyzed the user logs of pageviews, navigation patterns between photostreams. The authors used collaborative filtering method and built a stream transition graph to analyze common stream topic transitions to this end.

On retweet prediction, (Can, Oktay, and Manmatha 2013; Xu and Yang 2012; Petrovic, Osborne, and Lavrenko 2011) used classification-based approaches using tweet-based and author-based features. However, (Can, Oktay, and Manmatha 2013) took advantage of visual cues from images linked in the tweets, and (Xu and Yang 2012) employed social-based features in addition to tweet author-based features. Different from the other two works, (Xu and Yang 2012) performed the analysis from the perspective of individual users. (Petrovic, Osborne, and Lavrenko 2011) worked on retweet prediction of real-time tweeting with on-line learning algorithms and claimed that performance is dominated by social features, but that tweet features add a substantial boost. These studies showed that temporal features have a stronger effect on messages with low and medium volume of retweets compared to highly popular messages, and user activity features can further improve the performance marginally.

6 Conclusions

conclusion

7 Acknowledgments

We thank Kanna for her assistance in curating the topical hashtags in the Twitter data used in this study.

8 Copyright

References

- [Aiello et al. 2013] Aiello, L. M.; Petkos, G.; Martín, C. J.; Corney, D.; Papadopoulos, S.; Skraba, R.; Göker, A.; Kompatsiaris, I.; and Jaimes, A. 2013. Sensing trending topics in twitter. *IEEE Transactions on Multimedia* 15(6):1268–1282.

- [Albakour, Macdonald, and Ounis 2013] Albakour, M.-D.; Macdonald, C.; and Ounis, I. 2013. Identifying local events by using microblogs as social sensors. In *Proceedings of the 10th Conference on Open Research Areas in Information Retrieval, OAIR '13*.
- [Aramaki, Maskawa, and Morita 2011] Aramaki, E.; Maskawa, S.; and Morita, M. 2011. Twitter catches the flu: Detecting influenza epidemics using Twitter. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP '11*.
- [Becker, Naaman, and Gravano 2011] Becker, H.; Naaman, M.; and Gravano, L. 2011. Beyond trending topics: Real-world event identification on twitter. In *Proceedings of the Fifth International Conference on Weblogs and Social Media, Barcelona, Catalonia, Spain, July 17-21, 2011*.
- [Budak, Agrawal, and El Abbadi 2011] Budak, C.; Agrawal, D.; and El Abbadi, A. 2011. Structural trend analysis for online social networks. *PVLDB* 4(10):646–656.
- [Can, Oktay, and Manmatha 2013] Can, E. F.; Oktay, H.; and Manmatha, R. 2013. Predicting retweet count using visual cues. In *22nd ACM International Conference on Information and Knowledge Management, CIKM'13, San Francisco, CA, USA, October 27 - November 1, 2013*, 1481–1484.
- [Chen et al. 2012] Chen, K.; Chen, T.; Zheng, G.; Jin, O.; Yao, E.; and Yu, Y. 2012. Collaborative personalized tweet recommendation. In *Proceedings of the 35th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '12*, 661–670. New York, NY, USA: ACM.
- [Cheong and Cheong 2011] Cheong, F., and Cheong, C. 2011. Social media data mining: A social network analysis of tweets during the 2010-2011 australian floods. In *Pacific Asia Conference on Information Systems, PACIS 2011: Quality Research in Pacific Asia, Brisbane, Queensland, Australia, 7-11 July 2011*, 46.
- [Chiarandini et al. 2013] Chiarandini, L.; Grabowicz, P. A.; Trevisiol, M.; and Jaimes, A. 2013. Leveraging browsing patterns for topic discovery and photostream recommendation. In *Proceedings of the Seventh International Conference on Weblogs and Social Media, ICWSM 2013, Cambridge, Massachusetts, USA, July 8-11, 2013*.
- [Cui et al. 2012] Cui, A.; Zhang, M.; Liu, Y.; Ma, S.; and Zhang, K. 2012. Discover breaking events with popular hashtags in Twitter. In *21st ACM International Conference on Information and Knowledge Management, CIKM'12, Maui, HI, USA, October 29 - November 02, 2012*, 1794–1798.
- [Culotta 2010] Culotta, A. 2010. Towards detecting influenza epidemics by analyzing Twitter messages. In *Proceedings of the First Workshop on Social Media Analytics, SOMA '10*.
- [Diplaris et al. 2012] Diplaris, S.; Papadopoulos, S.; Kompatsiaris, I.; Göker, A.; MacFarlane, A.; Spangenberg, J.; Hacid, H.; Maknavicius, L.; and Klusch, M. 2012. Socialsensor: sensing user generated input for improved media discovery and experience. In *Proceedings of the 21st World Wide Web Conference, WWW 2012, Lyon, France, April 16-20, 2012 (Companion Volume)*, 243–246.
- [Fan et al. 2008] Fan, R.; Chang, K.; Hsieh, C.; Wang, X.; and Lin, C. 2008. LIBLINEAR: A library for large linear classification. *Journal of Machine Learning Research* 9:1871–1874.
- [Feld 1991] Feld, S. L. 1991. Why your friends have more friends than you do. *American Journal of Sociology* 1464–1477.
- [García-Herranz et al. 2012] García-Herranz, M.; Egido, E. M.; Cebrián, M.; Christakis, N. A.; and Fowler, J. H. 2012. Using friends as sensors to detect global-scale contagious outbreaks. *PloS one* abs/1211.6512.
- [Ishikawa et al. 2012] Ishikawa, S.; Arakawa, Y.; Tagashira, S.; and Fukuda, A. 2012. Hot topic detection in local areas using Twitter and wikipedia. In *ARCS Workshops (ARCS), 2012*, 1–5.
- [Krestel et al. 2015] Krestel, R.; Werkmeister, T.; Wiradarma, T. P.; and Kasneci, G. 2015. Tweet-recommender: Finding relevant tweets for news articles. In *Proceedings of the 24th International Conference on World Wide Web, WWW '15 Companion*, 53–54. Republic and Canton of Geneva, Switzerland: International World Wide Web Conferences Steering Committee.
- [Kryvasheyeyu et al. 2014] Kryvasheyeyu, Y.; Chen, H.; Moro, E.; Hentenryck, P. V.; and Cebrián, M. 2014. Performance of social network sensors during hurricane sandy. *PLoS one* abs/1402.2482.
- [Mathioudakis and Koudas 2010] Mathioudakis, M., and Koudas, N. 2010. Twittermonitor: trend detection over the Twitter stream. In *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2010, Indianapolis, Indiana, USA, 1155–1158*.
- [Nichols, Mahmud, and Drews 2012] Nichols, J.; Mahmud, J.; and Drews, C. 2012. Summarizing sporting events using Twitter. In *17th International Conference on Intelligent User Interfaces, IUI '12, Lisbon, Portugal, February 14-17, 2012*, 189–198.
- [O'Connor, Krieger, and Ahn 2010] O'Connor, B.; Krieger, M.; and Ahn, D. 2010. Tweetmotif: Exploratory search and topic summarization for Twitter. In *Proceedings of the Fourth International Conference on Weblogs and Social Media, ICWSM 2010, Washington, DC, USA, May 23-26, 2010*.
- [Petrović, Osborne, and Lavrenko 2010] Petrović, S.; Osborne, M.; and Lavrenko, V. 2010. Streaming first story detection with application to twitter. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics, HLT '10*, 181–189. Stroudsburg, PA, USA: Association for Computational Linguistics.
- [Petrovic, Osborne, and Lavrenko 2011] Petrovic, S.; Osborne, M.; and Lavrenko, V. 2011. Rt to win! predicting message propagation in Twitter. In *ICWSM*.
- [Phuvipadawat and Murata 2010] Phuvipadawat, S., and Murata, T. 2010. Breaking news detection and tracking in Twitter. In *Proceedings of the 2010 IEEE/WIC/ACM Inter-*

- national Conference on Web Intelligence and International Conference on Intelligent Agent Technology - Workshops, Toronto, Canada, August 31 - September 3, 2010*, 120–123.
- [Sadilek, Kautz, and Silenzio 2012] Sadilek, A.; Kautz, H. A.; and Silenzio, V. 2012. Modeling spread of disease from social interactions. In *Proceedings of the Sixth International Conference on Weblogs and Social Media, Dublin, Ireland, June 4-7, 2012*.
- [Sakaki et al. 2012] Sakaki, T.; Matsuo, Y.; Yanagihara, T.; Chandrasiri, N.; and Nawa, K. 2012. Real-time event extraction for driving information from social sensors. In *Cyber Technology in Automation, Control, and Intelligent Systems (CYBER), 2012 IEEE International Conference on*, 221–226.
- [Sakaki, Okazaki, and Matsuo 2013] Sakaki, T.; Okazaki, M.; and Matsuo, Y. 2013. Tweet analysis for real-time event detection and earthquake reporting system development. *Knowledge and Data Engineering, IEEE Transactions on* 25(4):919–931.
- [Starbird and Palen 2010] Starbird, K., and Palen, L. 2010. *Pass it on?: Retweeting in mass emergency*. International Community on Information Systems for Crisis Response and Management.
- [Vieweg et al. 2010] Vieweg, S.; Hughes, A. L.; Starbird, K.; and Palen, L. 2010. Microblogging during two natural hazards events: what Twitter may contribute to situational awareness. In *Proceedings of the 28th International Conference on Human Factors in Computing Systems, CHI 2010, Atlanta, Georgia, USA, April 10-15, 2010*, 1079–1088.
- [Wei et al. 2015] Wei, W.; Joseph, K.; Lo, W.; and Carley, K. M. 2015. A bayesian graphical model to discover latent events from twitter. In *Proceedings of the Ninth International Conference on Web and Social Media, ICWSM 2015, University of Oxford, Oxford, UK, May 26-29, 2015*, 503–512.
- [Weng and Lee 2011] Weng, J., and Lee, B. 2011. Event detection in Twitter. In *Proceedings of the Fifth International Conference on Weblogs and Social Media, Barcelona, Catalonia, Spain, July 17-21, 2011*.
- [Xu and Yang 2012] Xu, Z., and Yang, Q. 2012. Analyzing user retweet behavior on Twitter. In *International Conference on Advances in Social Networks Analysis and Mining, ASONAM 2012, Istanbul, Turkey, 26-29 August 2012*, 46–50.
- [Yamamoto and Satoh 2015] Yamamoto, S., and Satoh, T. 2015. Hierarchical estimation framework of multi-label classifying: A case of tweets classifying into real life aspects. In *Proceedings of the Ninth International Conference on Web and Social Media, ICWSM 2015, University of Oxford, Oxford, UK, May 26-29, 2015*, 523–532.
- [Yan, Lapata, and Li 2012] Yan, R.; Lapata, M.; and Li, X. 2012. Tweet recommendation with graph co-ranking. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers - Volume 1, ACL '12*, 516–525. Stroudsburg, PA, USA: Association for Computational Linguistics.
- [Zhao et al. 2011] Zhao, S.; Zhong, L.; Wickramasuriya, J.; and Vasudevan, V. 2011. Human as real-time sensors of social and physical events: A case study of Twitter and sports games. *Technical Report TR0620-2011, Rice University and Motorola Mobility* abs/1106.4300.