

# Learning Topical Social Sensors

**Authors**

Affiliations

## Abstract

Social media sources such as Twitter represent a massively distributed social sensor of a rich underlying topic space ranging from social and political events to entertainment and sports news. However, given the continual evolution of social media content, querying for content from individual users or containing certain keywords or hashtags is often insufficient to retrieve the vast range of topical content available. To automate generic social sensor construction, we train a supervised learner to identify topical content over a large feature space given a small set of seed hashtags. On a corpus of approximately (choose one: 1 billion English Tweets / 40 TB of data) collected from the Twitter streaming API during 2013 and 2014 and learning for 10 diverse topics 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 (including content with no hashtags) 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 performance.

## 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 “Social Sensor” research<sup>1</sup>), automatically learning high-precision sensors (i.e., ranking and retrieval meth-

ods) 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, manual labeling is simply too time-consuming for many users and crowdsourced labels are both costly and prone to misinterpretation of users’ information needs. Fortunately, 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 (choose one: 1 billion English Tweets / 40 TB of data) 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.

Tweets					
Feature	Max	Avg	Median	Max entity	Count
From	10,196	8.67	2	running_status	95,547,198
Hashtag	1,653,159	13.91	1	#retweet	11,183,410
Mention					411,341,569
Location					58,601
Term	2024529	7,450.58	323	the	20,234,729
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	6,616.65	305	the	
Hashtags					
From	18,167	1.62	0	daily_astrodata	

Table 1: Feature Statistics

## Dataset Statistics

We crawled Twitter data using Twitter Streaming API for two years spanning 2013 and 2014 years. This type of crawling provides us with a very sparse set of data, roughly 1% of all tweets<sup>2</sup>. 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. These uni-grams are later cleaned to remove *Terms* with no meaning (total number of *Terms* before cleaning was 20,234,729)

Table 1 provides more detail statistics about each feature. For each feature, we reported the count of the feature in our dataset, in addition to maximum, average, median counts of each feature across the tweets. Lower part of the table provides these counts across user dimension meaning that for example a hashtag has been used in average by 10.08 users. Last part of the table shows the statistics for the hashtag usage of our users e.g., users have used 2 hashtags in average.

Figure ?? shows details of number of tweets per month and figure ?? shows the power law plots of tweet counts and hashtag counts for users. We chose 10 topics for our experiments. Tweets are temporally divided over 2 years to provide train and test sets. Table 2 provides samples of training hashtags and number of train hashtags, test hashtags, topical tweets for each topic. Some topics such as *HumanCausedDisaster* and *Soccor* 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.

Figure ?? shows distribution of tweets across different location in U.S. and international locations overall and for each topic(?).

<sup>2</sup>[http://allthingsd.com/20101110/twitter-firehose-too-intense-take-a-sip-from-the-garden-hose-or-sample-the-spritzer—](http://allthingsd.com/20101110/twitter-firehose-too-intense-take-a-sip-from-the-garden-hose-or-sample-the-spritzer/)

## Experimental Methodology

How we curated hashtags: need to make up good story here. Inner-annotator agreement of 3/4.

Train/validation/test split date selection – temporally 5,.1,.4

Feature selection: threshold per feature 159 and 50 (just explain rationale for lower hashtag and location thresholds).

Formal notation, how do we train/test and tune hyperparameters for a generic classifier.

## Classification Algorithms

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

All above over 1,000,000 features, \*same\* training data for all algorithms.

Not breaking down by feature type yet – that’s for the feature analysis section.

## Analysis

- Table of rows:alg, cols: MAP, P@k (k in 10,100,1000) with stderrs over all topics

- Could do a bar graph (below) each for MAP, P@100 with topics as major columns and algs as neighboring bars

- Anecdotal results for each topic – point out deficiency in our labels (a good thing, we generalized well from small hashtag set), manual evaluation of relevance for top-100 for best algorithm?

## 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:

$$I(t_i, f_k) = \sum_{t_i \in \{true, false\}} \sum_{f_k \in \{true, false\}} p(f_k, t_i) \log \left( \frac{p(f_k, t_i)}{p(f_k)p(t_i)} \right) \quad (1)$$

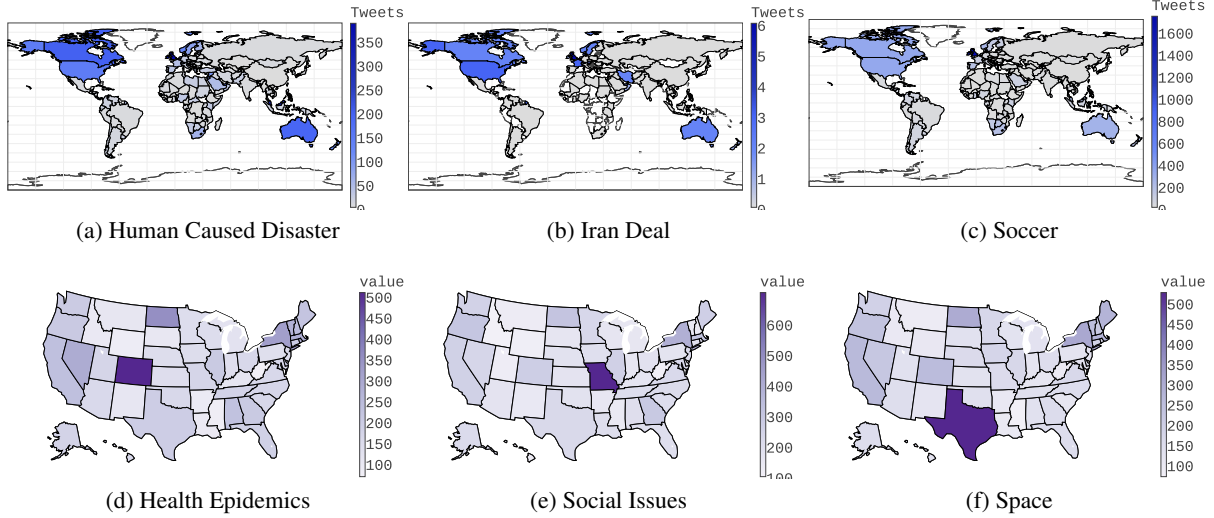


Figure 1: Choropleths of top: International map and bottom: U.S. map

Topics/Top10	NaturalDisaster	Epidemics	IranDeal	SocialIssues	LBGT	HumanCausedDisaster	CelebrityDeath	Space	Tennis	Soccer
#TrainHashtags	31	52	12	31	29	49	28	98	58	126
#TestHashtags	18	33	5	19	17	29	16	63	36	81
#TotalTopicalTweets	42,987	210,217	8,762	230,058	282,527	408,304	163,890	239,719	55,053	860,389
Sample Train Hashtags	#earthquake	#ebola	#irandeal	#policebrutality	#loveislove	#gazaundersattack	#robinwilliams	#asteroids	#usopenchampion	#worldcup
	#storm	#virus	#iranfreedom	#michaelbrown	#gaypride	#childrenofsyria	#ripmandela	#astronauts	#novakdjokovic	#lovesoccer
	#tsunami	#vaccine	#irantalk	#justice4all	#uniteblue	#iraqwar	#ripjoanrivers	#satellite	#wimbledon	#fifa
	#abloods	#chickenpox	#rouhani	#freetheweed	#homo	#bombthreat	#mandela	#spacecraft	#womenstennis	#realmadrid
	#hurricanekatrina	#theplague	#nuclearpower	#newnjgunlaw	#gaymarriage	#isis	#paulwalker	#telescope	#tennisnews	#beckham

Table 2: Test/Train Hashtag samples and statistics

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

First, we provide mutual information values for each feature across different topics shown by boxplots in figure ??, and average values of mutual information for each feature vs different topics shown in table 2. The last column in table 2 shows average mutual information for the feature with the standard error range provided. We make a few observations from the analysis of Table 2:

- Term features provide more information for all of the topics on average which shows the importance of uni-grams when it comes to selection of topical tweet.
- From and mention features are the least informative features for all of the topics.
- 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 HumanCausedDisaster, LBGT, and Soccer topics.
- Sorting features based on their average mean value across different topics results in the following order:

1. Term
2. Location
3. Hashtag

4. Mention
5. From

	Tennis	Space	Soccer	IranDeal	HumanDisaster	CelebrityDeath	SocialIssues	NaturalDisaster	Epidemics	LBGT	Mean
Mention	0.06	0.23	0.61	0.02	0.46	0.11	0.30	0.05	0.19	0.34	0.24
Hashtag	0.84	3.40	9.54	0.28	5.63	1.84	2.88	0.66	3.54	4.93	3.35
From	0.05	0.19	0.53	0.01	0.29	0.12	0.16	0.04	0.17	0.22	0.18
Location	5.24	16.05	42.94	1.04	31.82	6.03	17.82	4.03	15.17	32.57	17.08
Term	17.10	45.49	144.27	4.48	128.82	32.63	68.34	13.22	44.67	73.37	57.24

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

It is important to note that due to very large amount of Term features, they were cleaned based on their frequency (having at least frequency value of 100).

**For each feature type, do any attributes correlate with importance?** In order to give a better sense of what features are better for each topic, we provided top-5 features

Method	Metric	Tennis	Space	Soccer	IranDeal	HumanCausedDisaster	CelebrityDeath	SocialIssues	NaturalDisaster	Epidemics	LGBT	Mean±Std
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 3: Different learning methods results on topics with hyper-parameter tuning based on MAP

<b>Tennis</b>	<b>Space</b>
✓ rt @esptennis: 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
✓ rt @esptennis: shock city, darcis drops rafa in straight sets. first time nadal loses in first rd of a. major in career. #espnwimbledon #w	✓ voting mars @30secondstomars @jaredleto @shannonleto @tomofromearth xobest group http://t.co/dlsozvjnf
✓ rt @esptennis: shock city, darcis drops rafa in straight sets. first time nadal loses in first rd of a. major in career. #espnwimbledon #w	✓ rt @jaredleto.com: show everyone how much you are proud of @30secondstomars #mthottest 30 seconds to mars http://t.co/byxni467
✓ rt @esptennis: shock city, darcis drops rafa in straight sets. first time nadal loses in first rd of a. major in career. #espnwimbledon #w	✓ rt @30secondstomars: missed the big news? mars touring with @linkinpark + special guests @afi this summer! http://t.co/3e5m0pwwd
✓ rt @esptennis: shock city, darcis drops rafa in straight sets. first time nadal loses in first rd of a. major in career. #espnwimbledon #w	✓ rt @30secondstomars: to the right! to the left, we will fight to the death! #intothewildonmyrt with mars, starting weekly, nov 30 _tn
<b>Soccer</b>	<b>IranDeal</b>
✓ rt @tomn_dogg: #thingstodobeforeearthends spend all my money.	✓ rt @iran_policy: @vidalquadrax: @isjcommittee has investigated 10 major subjects of irans controversial #nuclear program #iranalksvienna
✓ rt @mancivilian: nice performance	✓ rt @iran_policy: @vidalquadrax: @isjcommittee has investigated 10 major subjects of irans controversial #nuclear program #iranalksvienna
✓ rt @indykalla: pedokali: "It's see what happens in the winter. the fact is that i'm not happy with it, that's clear." @arsenal	✓ rt @negamortazavi: thank you @hasanmohammadi for reweeting. 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/sdrakqkly #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. #iranviolations http://t.co/1gdxh1znu
<b>HumanCausedDisaster</b>	<b>CelebrityDeath</b>
✓ rt @baselyrian: there've been peaceful people in #homs not terrorists! #assad, enemy of #humanity destroyed it. #eyehomhs #withsyria http://	✓ rt @savubona.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! #peakup4syrianchildren 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& his country, he can res
✓ exclusive: us formally requested #un investigation: russia pressured #assad to no avail, chain of evidence proof hard http://t.co/5602rvdfw	✓ 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
✓ #savealeppo from #assadwarcrimes #savealeppo from #civilians - targeted shelling of #assad regime #syria #aleppo http://t.co/3dfrh0pnl	✓ rt @jacques #kallis: a phenomenal cricketing giant of all time - #cricket #history #southafrica http://t.co/ms5pmwoag9
✓ rt @maine_rights: why does the #un allow this to continue? rt @tintin1957 help raise awareness of the suffering in #syriawarcrimes http://t	✓ rt @sudehi304: south africa has the most beautiful babies...so diverse, so unique...so god!! lol #durban #southafrica
<b>SocialIssues</b>	<b>NaturalDisaster</b>
✓ the us doesn't actually borrow is the thing, i believe in a creationist theory of the us dollar @usanationdebt @nationaldebt	✓ rt @jacknecolcute: @lunaticrex @fingersmalloy @toddincannon @theanonliberal anthony kennedy just wrote opinion granting legal protection to cupcake kiplers
✓ rt @2anow: according to @njenatetpes women's rights do not include this poor nj mother's right to defend herself http://t.co/czbslnqk66	✓ rt @toddincannon: your personal account, your interest, separate from your business.
✓ rt @2anow: confiscation? how many carry permits are in the senate and assembly? give us ours or turn them in. @senatorforetas @lougreenw	✓ why would you report someone as spam if he is not spam? @illygirlbrea @toddincannon
✓ rt @2anow: vote with your wallet against #guncontrolforever city enterprises does not support the #2a http://t.co/tpk0k3berrm9j2as #cot	✓ rt @3b_arch3r: @toddincannon thanks for your if having the female realbrother. between them is 600 lbs, 104 iq points, and a lot of hate.
✓ rt @2anow: @mensdemand @jtimes3 they dont have a plan for that, which is why they should never be allowed to take our guns	✓ rt @toddincannon who us dick tickle.
<b>Epidemics</b>	<b>LGBT</b>
✓ rt @who: fourteen of the susp. & conf. ebola cases in #conakry, #guinea, are health care workers, of which 11 died #askubola	
✓ who who who can afford also been cover in government health insurance [with universal health coverage]	
✓ #ebolabreak this health crisis. unparalleled in modern times. @who dir. gzyward - requires \$1 billion to stem http://t.co/fqgdyb83d	
✓ rt @medsin: @who are conducting a survey on the social determinants of health in medical teaching. fill the survey in at https://t.co/g559y	
✓ augmentation vertigineuse de 57,4% en 1 an des actes islamophobes en france, dit le collectif contre l'islamophobie http://t.co/2qhgocg55	

Table 4: Top Tweets for each topic based on MAP tuned results

for each topic in table 5. 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.

**scatterplots of feature MI – the absolute last thing we do (density plots??!!) \*\*which plots below, and for which topics? Could pick out most useful features for topics in part (a)(i) and just show selected scatter plots below for these feature types. from, mention MIs vs. followers, favorites, friends, hashtags, tweets hashtag MI vs. #tweets, #users location MI vs. #users term MI vs. #tweets**

## 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-

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 second 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

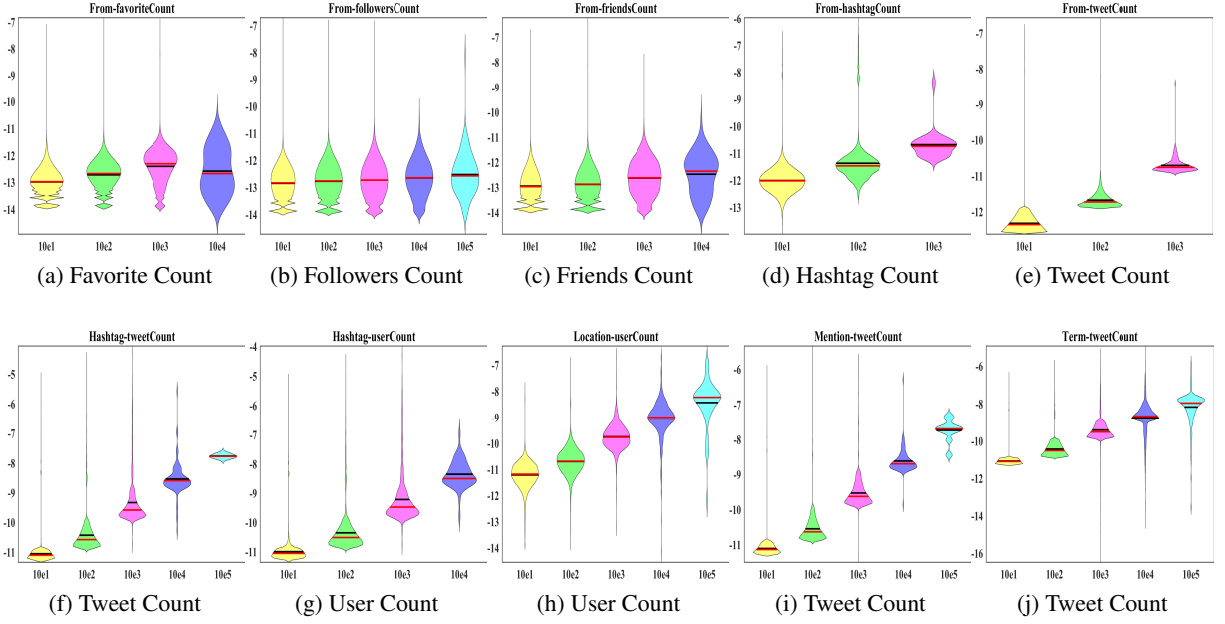


Figure 3: ViolinPlots for feature attributes counts vs. MI. Top row shows attributes {favoriteCount, followerCount, friendCount, hashtagCount, tweetCount} for *From* feature. Bottom row shows attributes tweetCount and/or userCount for *Hashtag*, *Location*, *Mention*, and *Term* features.

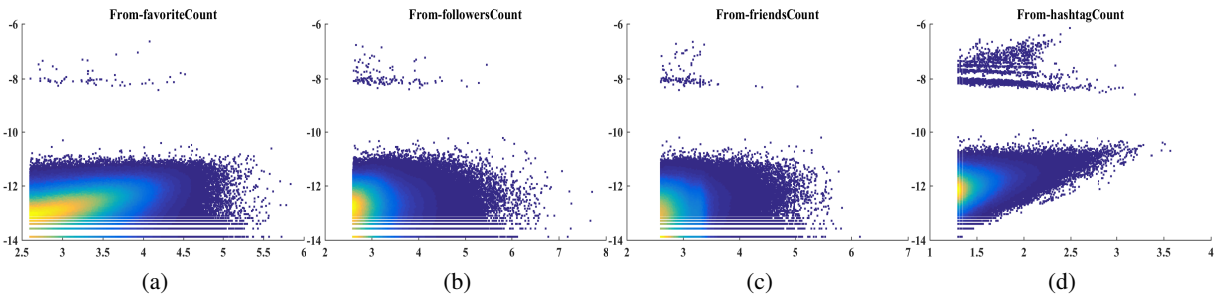


Figure 4: DensityPlots for feature attributes counts vs. MI. (a-d) show attributes {favoriteCount, followerCount, friendCount, hashtagCount} for *From* feature

Topics/Top10	NaturalDisaster	Epidemics	IranDeal	SocialIssues	LGBT	HumanCausedDisaster	CelebrityDeath	Space	Tennis	Soccer
From	earthquake_wo	changedecopine	mazandara	nsingerdebt	eph4_15	ydumozyf	nmandelaquotes	daily_astrodata	tracktennisnews	losangelessrh
From	earthalerts	drdaveanddee	hhadi119	debtadvisoruk	mgdauber	syriatweeten	boiknox	freesolarleads	tennis_result	shoetale
From	seelites	joimmentotetwk	140iran	debt_protect	stevendickinson	tintin1957	jacanews	houston_jobs	i_roger_federer	sport_agent
From	globalfloodnews	followebola	setarehgan	negativeequityf	lileensvf1	sirajsol	ewnreporter	star_wars_gifts	tennislessonnow	books_you_want
From	gcmcdrought	localnursejobs	akgharshabaneh	dolphin_ls	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	ripcoreymonteith	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	foxtamedia	4freedominiran	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	theanonmessage	barackobama	palestinianism	tweetlikegiris	30secondstomars	espn tennis	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	acrux	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 5: Top 5 features for each topic based on Mutual Information

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)<sup>3</sup>). (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

<sup>3</sup>On average, most people have fewer friends than their friends have

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 people. 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 online 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.

## Conclusions

conclusion

## Acknowledgments

acknowledgements

## 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.
- [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 International 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.