

Category-Level Transfer Learning from Knowledge Base to Microblog Stream for Accurate Event Detection

Weijing Huang, Tengjiao Wang, Wei Chen, Yazhou Wang

School of Electronics Engineering and Computer Science, Peking University

@DASFAA 2017, Suzhou, China

Motivation

Many Web applications need the **accurate event detection** technique on microblog stream, including:

- ① public opinion analysis [Chen, SIGIR 2013]
- ② public security [Li, ICDE 2012], [Imran, WWW 2014]
- ③ disaster response [Sakaki, WWW 2010]
- ④ breaking news report¹

But detecting events on twitter stream accurately is still challenging.

¹<http://www.theverge.com/2016/12/1/13804542/reuters-algorithm-breaking-news-twitter>

Challenges (1/2)

According to [Huang, WWW 2016], the challenges include,

- ① fast changing
- ② high noise
- ③ short length

And, we found another key factor,

- ① Small events with fewer tweets → Hard to trade off between precision and recall.

Challenges (2/2)

Exploratory study on the *Edinburgh twitter corpus*: 11/27 events contain less than 50 tweets.

Table: Statistics of labeled events.

Event	Date	Event Size
S&P downgrade US credit rating	05/08/2011	656
Atlantis shuttle lands	21/07/2011	595
US increases debt ceiling	25/07/2011	485
Plane with Russian hockey team Lokomotiv crashes	07/09/2011	286
Amy Winehouse dies	23/07/2011	283
Gunman opens fire in youth camp in Norway	23/07/2011	260
Earthquake in Virginia	24/08/2011	246
First victim of London riots dies	09/08/2011	174
Explosion in French nuclear plant in Marcoule	12/09/2011	135
Google announces plans to bury Motorola Mobility	15/08/2011	127
NASA announces there might be water on Mars	04/08/2011	124
Car bomb explodes in Oslo, Norway	22/07/2011	114
...
Indian and Bangladesh sign a border pact	06/09/2011	25
Flight 4896 crash	13/07/2011	21
First artificial organ transplant	12/07/2011	18
three men die in riots in england	10/08/2011	16
rebels capture interational tripoli airport	21/08/2011	13

11
Small
events
with
fewer
tweets

How about existing methods? (1/2)

Event detection methods *without extra information*, such as

- ① clustering articles
 - LSH[Petrovic, NAACL 2010]
 - need to set threshold to determine whether new article represents a new event.
- ② analyzing word frequencies
 - EDCoW[Weng, ICWSM 2011]
 - treat the word as the basic unit in analysis, without regarding polysemy words (words have different meanings, e.g. “apple”)
- ③ finding bursty topics via topic modeling
 - TimeUserLDA[Diao, ACL 2012], BurstyBTM[Yan, AAAI 2015]
 - detects the “large” events but may ignore the “small” ones.

How about existing methods? (2/2)

Event detection methods *by leveraging extra information*

- ① typical one: Twevent[Li, CIKM 2012]
 - divides the tweet into segments according to the Microsoft Web N-Gram service and Wikipedia
 - detects the bursty segments and cluster these segments into candidate events
 - still has to trade off between precision and recall
 - e.g., the bursty segment in the not-so-popular event “*first artificial organ transplant*” is missed

An Example

Much easier to detect the event on the time series of word “hood” related to *Military*, **without adjusting the threshold**.

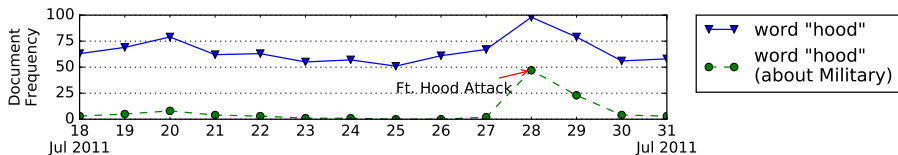


Figure: The comparison of the time series between the raw word *hood* and the *Military* related word *hood*, computed on the *Edinburgh twitter corpus*. Refer the event to https://en.wikipedia.org/wiki/Fort_Hood#2011_attack_plot.



An infant wearing a hood.

Fort Hood

From Wikipedia, the free encyclopedia

Fort Hood is a U.S. military post located in Killeen, Texas. The post is named after Confederate General John Bell Hood. It is

The insights on the example

Knowledge Base

- Well organized
- Constructed elaborately
- full of rich information

Microblog stream

- Short length
- Fast changing
- High noise

The benefit of **enriching the semantics** and **filtering out noise** by Knowledge Base for microblogs is attractive.

But it's expensive to retrieve every word of tweets in the Knowledge Base.

Our solution

TRANSDetector: a novel category-level transfer learning method

- transfer KB's **category-level info** into microblog stream
- balance the performance and the cost of leveraging knowledge base for event detection

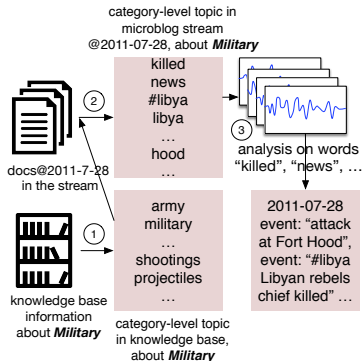


Figure: TRANSDetector's processing flow, in 3 phases.

TRANSDetector: Phase 1 (Extracting Category-Level Topics in KB) (1/3)

There is a three fold hierarchical structure in Knowledge Base.

- ① **Taxonomy Graph** $G^{(0)}$. The directed edges in $G^{(0)}$ represent the *class*→*subclass* relations in KB.
 - e.g., *Main topic classifications* → *Society*
- ② **Category-Page Bipartite Graph** $G^{(1)}$. The directed edges in $G^{(1)}$ represent the *class*→*instance* relations in KB.
 - e.g., *Military* → *page: Armed Forces*
- ③ **Page-Content Map** $G^{(2)}$. For a specific Wikipedia dumps version, the edges *page* → *content* in $G^{(2)}$ define a one-to-one mapping.
 - e.g., *page: Armed Forces* → *content(20170325version)*: “The armed forces of a country are its government-sponsored defense, fighting forces, and organizations...”

TRANSDetector: Phase 1 (Extracting Category-Level Topics in KB) (2/3)

Algorithm 1: Extraction of Category-Level Topics in Knowledge Base

Input: Taxonomy's Graph $G^{(0)}$, Category-Page Bipartite Graph $G^{(1)}$, Page-Content Bipartite Graph $G^{(2)}$, topic related category node c

Output: c 's category-level topic in knowledge base h_c

```
1  $Pages(c) \leftarrow \emptyset, h_c \leftarrow \emptyset$ 
2  $DAG\ G^{(0)'} \leftarrow$  Remove Cycles of  $G^{(0)}$  by nodes' HITS-PageRank scores.
3  $SuccessorNodes(c) \leftarrow$  Breadth-first-traverse( $G^{(0)'}$ ,  $c$ )
4 for  $node \in SuccessorNodes(c)$  do
5    $Pages(c) \leftarrow Pages(c) \cup G^{(1)}.neighbours(node)$ 
6 Word frequency table  $n(c, .) \leftarrow$  do word count on the text contents of  $Pages(c)$ 
7 Word frequency table  $n(All, .) \leftarrow$  do word count on the text contents of all pages in  $G^{(2)}$ .
8 for  $word\ w\ in\ WordFrequencyTable(All).keys()$  do
9    $chi(c, w) \leftarrow$   $w$ 's chi-square statistics on  $WordFrequencyTable(c)$  and  $WordFrequencyTable(All)$ .
10   $h_{c,w} \leftarrow chi(c, w)$ 
11 return  $h_c$ 
```

TRANSDetector: Phase 1 (Extracting Category-Level Topics in KB) (3/3)

Taking the category *Military* as an example, we extract *Military*'s category-level topic h_{Military} .

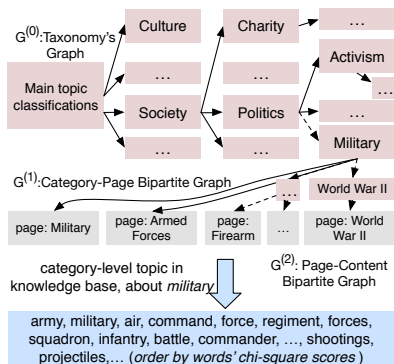


Figure: Extracting Category-Level Topics in Knowledge Base via its three fold hierarchical structure, taking *Military* as an example.

TRANSDetector: Phase 2 (Transferring Category-Level Info into Microblog Stream) (1/2)

Transfer KB's Category-Level Topics $\{\mathbf{h}_c\}_{c=1}^{K_{KB}}$ into microblogs stream: CTrans-LDA.

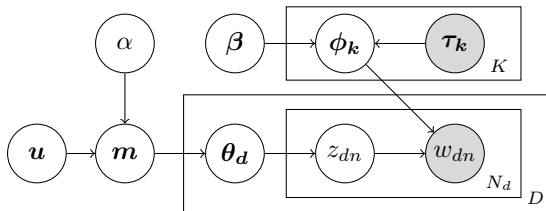


Figure: Diagram of CTrans-LDA.

In CTrans-LDA, $\{\mathbf{h}_c\}_{c=1}^{K_{KB}}$ is used as prior information:

$$\tau_{kv} = \begin{cases} \lambda \frac{h_{kv}}{\sum_{v \in S_k} h_{kv}}, v \in S_k \text{ and } k \leq K_{KB} \\ 0, v \notin S_k \text{ or } k > K_{KB} \end{cases} \quad (1)$$

TRANSDetector: Phase 2 (Transferring Category-Level Info into Microblog Stream) (2/2)

We use Gibbs Sampling for solving CTrans-LDA.

- The initialization probability $\hat{q}_{k|v}$ makes sure that the learned topics are aligned to the pre-defined category-level topic.

$$\hat{q}_{k|v} = \begin{cases} \frac{\tau_{kv}}{\sum_{k=1}^K \tau_{kv}}, \sum_k \tau_{kv} > 0 & (a) \\ 0, \sum_k \tau_{kv} = 0 \text{ and } k \leq K_{KB} & (b) \\ 1/(K - K_{KB}), \sum_k \tau_{kv} = 0 \text{ and } k > K_{KB} & (c) \end{cases} \quad (2)$$

- Conditional probability in gibbs sampling:

$$p(z_{dn} = k | \cdot) \propto (n_{dk}^{(d)} + \alpha m_k)(n_{kv}^{(w)} + \tau_{kv} + \beta) / (n_{k,\cdot}^{(w)} + \tau_{k,\cdot} + V\beta).$$

TRANSDetector: Phase 3 (Detecting Events on Category-Level Word Time Series) (1/2)

After transfer learning, we conduct analysis on category-level word time series, and detect events in microblog stream.

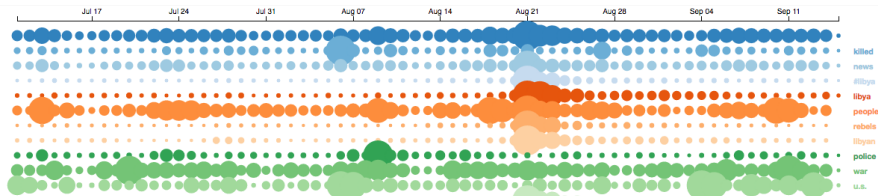


Figure: Visualizing Category-Level Word Time Series in Microblog Stream on Edinburgh Twitter Corpus (20110711-20110915), taking *Military* as an example.

TRANSDetector: Phase 3 (Detecting Events on Category-Level Word Time Series) (2/2)

With richer semantics, and fewer noise, we detect events more accurately, in all the following granularities.(see more technique details in our paper)

- ① events' candidate words
 - e.g., *Ft.*, *Hood*, *attack*.
- ② event phrases
 - e.g., *Ft. Hood attack*.
- ③ events' representative microblogs
 - e.g., *Possible Ft. Hood Attack Thwarted <http://t.co/BSJ33hk>*.

Experiment Settings

Dataset

- Knowledge Base. Wikipedia dumps²³
- Microblog Stream. Edinburgh twitter corpus⁴

Baseline Methods

- Twevent, BurstyBTM, LSH, EDCoW, and TimeUserLDA

Ground Truth

- Benchmark1 is labeled events in previous study.
- Benchmark2 is our manually checked events based on Twevent, BurstyBTM, LSH, EDCoW, TimeUserLDA and TRANSDETECTOR.

²<https://dumps.wikimedia.org/enwiki/latest/enwiki-latest-categorylinks.sql.gz>

³<https://dumps.wikimedia.org/enwiki/latest/enwiki-latest-pages-articles.xml.bz2>

⁴http://demeter.inf.ed.ac.uk/cross/docs/fsd_corpus.tar.gz

Experimental Results

Evaluation on Category-Level Topics in Knowledge Base. On *Aviation* topic, semantic coherence is much better than LightLDA (same as LDA) in terms of NPMI[Roder, WSDM 2015]⁵.

Table: The comparison on the topic coherence(NPMI) between our method and LightLDA, taking *Aviation* as an example. (NPMI is computed on a group of ten words. \sim stands for the top five words.)

Category-Level Topics extracted from Wikipedia by TRANSDETECTOR				Topics Learned from Wikipedia by LightLDA			
GID	#words*	words	NPMI	GID	#words*	words	NPMI
-	1-5	aircraft air airport flight airline	-	-	1-5	engine aircraft car air power	-
0	1-5, 6-10	\sim , airlines aviation flying pilot squadron	0.113	0	1-5, 6-10	\sim , design flight model production speed	0.112
1	1-5, 11-15	\sim , flights pilots raf airways fighter	0.155	1	1-5, 11-15	\sim , system vehicle cars engines mm	0.062
2	1-5, 16-20	\sim , boeing runway force crashed flew	0.092	2	1-5, 16-20	\sim , fuel vehicles designed models type	0.072
3	1-5, 21-25	\sim , airfield landing passengers plane aerial	0.179	3	1-5, 21-25	\sim , version front produced rear electric	0.035
4	1-5, 26-30	\sim , bomber radar wing bombers crash	0.137	4	1-5, 26-30	\sim , space control motor standard development	0.085
5	1-5, 31-35	\sim , airbus airports operations jet helicopter	0.189	5	1-5, 31-35	\sim , film range light using available	-0.002
6	1-5, 36-40	\sim , squadrons base flown havilland crew	0.088	6	1-5, 36-40	\sim , wing powered wheel weight launch	0.087
7	1-5, 41-45	\sim , combat luftwaffe aerodrome carrier fokker	0.159	7	1-5, 41-45	\sim , developed low test ford cylinder	0.007
8	1-5, 46-50	\sim , planes fly engine takeoff fleet	0.186	8	1-5, 46-50	\sim , equipment side pilot hp aviation	0.091
9	1-5, 51-55	\sim , fuselage helicopters aviator naval aero	0.157	9	1-5, 51-55	\sim , systems us sold body drive	-0.051
10	1-5, 56-60	\sim , glider command training balloon faa	0.166	10	1-5, 56-60	\sim , gear introduced class safety seat	0.069
...
18	1-5, 96-100	\sim , scheduled carriers military curtiss biplane	0.131	18	1-5, 96-100	\sim , transmission special replaced limited different	0.059
19	1-5, 101-105	\sim , accident engines iaf albatross rcaf	0.068	19	1-5, 101-105	\sim , features machine nuclear even unit	0.011

⁵<https://github.com/AKSW/Palmetto>

Experimental Results

Evaluation on Category-Level Topics in Knowledge Base. On more topics, semantic coherence is much better than LightLDA (same as LDA) in terms of NPMI[Roder, WSDM 2015].

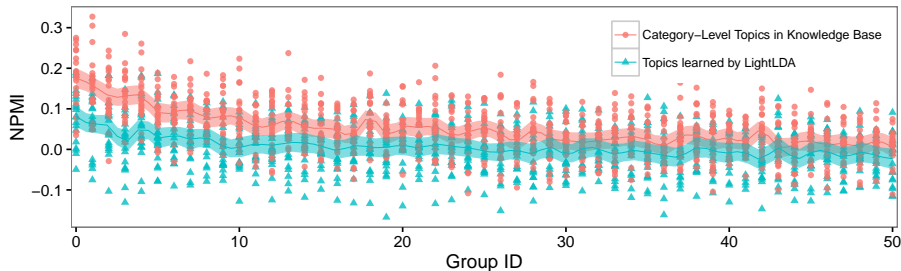


Figure: More topics are compared at the NPMI metrics between our method and LightLDA

Experimental Results

Effectiveness of transferring category-level topics into the microblog stream, and finding more new words in the stream which is not stored in the knowledge base.

Table: Category-Level Topics extracted from knowledge base and the corresponding topics on microblog stream learned from CTrans-LDA. The words in **bold** font are newly learned on the microblog stream by the transfer learning.

<i>Aviation</i>		<i>Health</i>		<i>Middle East</i>		<i>Military</i>		<i>Mobile Phones</i>	
Knowledge Base	Microblog Stream	Knowledge Base	Microblog Stream	Knowledge Base	Microblog Stream	Knowledge Base	Microblog Stream	Knowledge Base	Microblog Stream
aircraft	air	health	weight	al	#syria	army	killed	android	iphone
air	plane	patients	loss	israel	#bahrain	military	news	mobile	apple
airport	flight	medical	diet	iran	people	air	#libya	nokia	android
flight	time	disease	health	arab	israel	command	libya	ios	app
airline	airlines	treatment	cancer	israeli	police	force	rebels	phone	ipad
airlines	news	hospital	lose	egypt	#libya	regiment	people	samsung	samsung
aviation	boat	patient	fat	egyptian	#egypt	forces	police	game	mobile
flying	airport	clinical	tips	ibn	news	squadron	war	app	blackberry
pilot	force	symptoms	treatment	jerusalem	#israel	infantry	libyan	iphone	tablet
squadron	fly	cancer	body	syria	world	battle	attack	htc	apps

Experimental Results

Effectiveness of event detection:

- 1 TRANSDETECTOR performs better in terms of the precision and the recall.
- 2 only sacrificing in the DERate slightly because an event could be grouped into multiple categories.
 - the event “*S&P downgrade US credit rating*”, related to *politics* and *financial* simultaneously.

Table: Overall Performance on Event Detection

Method	Number of Events to be Evaluated	Recall@ Benchmark1	Precision@ Benchmark2	Recall@ Benchmark2	F@ Benchmark2	DERate ^a (Duplicate Event Rate)@ Benchmark2
LSH	500	0.704	0.788	0.651	0.713	0.348
TimeUserLDA	100	0.370	0.790	0.177	0.289	0.114
Twevent	375	0.741	0.808	0.658	0.725	0.142
EDCoW	349	0.556	0.748	0.511	0.607	0.226
BurstyBTM	200	0.667	0.825	0.384	0.497	0.079
TRANSDETECTOR	457	0.889	0.912	0.876	0.894	0.170

^a DERate = (the number of duplicate events) / (the total number of detected realistic events)

Experimental Results

To understand why TRANSDETECTOR performs better.

Show the relationship between the recall and the event size.

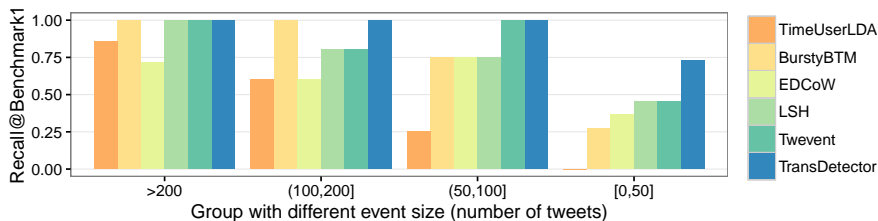


Figure: The relation between the recall and the event size

Experimental Results

To understand why TRANSDetector performs better.
Show the relationship between the recall and the event size, taking the *military*-related events as an example.

Table: Events about *military* detected by systems between 2011-07-22 and 2011-07-28

Date	Event key words	Representative event tweet	Number of event tweet	Methods ^a					
				L	TU	TW	E	B	TD
7/22/11	Norway, Oslo, attacks, bombing	Terror Attacks Devastate Norway: A bomb ripped through government offices in Oslo and a gunman... http://dlvr.it/cLbk8	557	✓	✓	✓	✓	✓	✓
7/23/11	Gunman, rink	Gunman Kills Self, 5 Others at Texas Roller Rink http://dlvr.it/cLcTH	43	-	-	✓	✓	-	✓
7/26/11	Kandahar, mayor, suicide, attack	TELEGRAPH]: Kandahar mayor killed by Afghan suicide bomber: The mayor of Kandahar, the biggest city in south _	47	✓	-	✓	✓	-	✓
7/28/11	Ft., Hood, attack	Possible Ft. Hood Attack Thwarted http://t.co/BSJ33hk	52	-	-	-	-	-	✓
7/28/11	Libyan, rebel, gunned	Libyan rebel chief gunned down in Benghazi http://sns.mx/prfvy1	44	-	-	-	-	-	✓

^a L=LSH, TU=TimeUserLDA, TW=Twevent, E=EDCoW, B=BurstyBTM, TD=TRANSDetector.

Conclusions

- Knowledge base is constructed elaborately and contains rich information, which can benefit the not-well-organized microblog stream.
- We propose TRANSDETECTOR method, and
 - use category-level topic in knowledge base as the prior knowledge,
 - transfer abundant knowledge from knowledge base into microblog stream
 - enrich the semantics of microblogs and further enhances the accuracy of microblogs event detection

Thanks!

Q&A

⁵This slide and more data are available at <http://q-r.to/bajx8I>