

An Empirical Study on Uncertainty Identification in Social Media Context

**Zhongyu Wei¹, Junwen Chen¹, Wei Gao², Binyang Li¹
Lanjuan Zhou¹, Yulan He³, and Kam-Fai Wong¹**

¹The Chinese University of Hong Kong

²Qatar Computing Research Institute, Qatar Foundation, Doha, Qatar

³School of Engineering, Applied Science, Aston University, Birmingham, UK

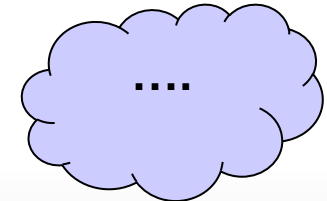
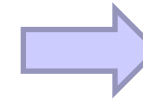
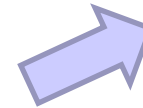
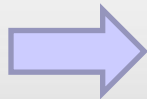
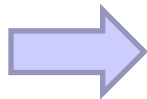
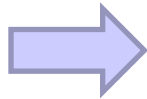
August 5th, 2013 at Sofia, Bulgaria

The 51st Annual Meeting of the Association for Computational Linguistics

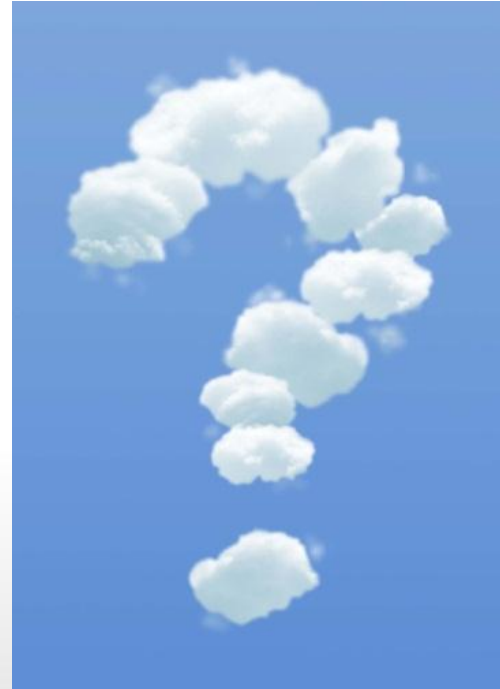
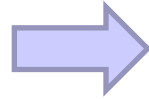
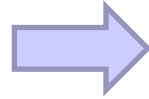
The Chinese University of Hong Kong
Department of Systems Engineering & Engineering Management



Background



Background



Factuality



Uncertainty

- “Uncertainty” can be interpreted as lack of information: the receiver of the information (i.e., the hearer or the reader) cannot be certain about some pieces of information”.



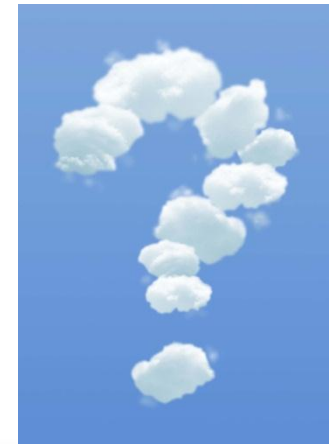
Uncertainty

- Related work
 - Binary uncertainty classification on formal text.
 - CoNLL shared task 2010
 - Existing uncertainty corpus.
 - *Factbank (Newswires)*
 - *BioScope (Biology paper)*
 - *Wikipedia Weasels (Wikipedia article)*



Motivation

- 2011 London Riots dataset
 - 18.9% of 326,747 tweets contain uncertainty keyword
- Rare work on social media
 - Uncertainty identification is domain dependent.
 - No corpus available in social media context.



*Probably
Possibly
Maybe*

...

...



Contribution

- We propose a variant of classification scheme for uncertainty identification in social media context.
- We construct the first uncertainty dataset in social media context.
- We perform uncertainty identification experiments and explore effectiveness of different types of features.



Traditional Classification*

- **Epistemic:**

- On the basis of our world knowledge we cannot decide at the moment whether the statement is true or false.
 - **Possible:** *It **may be** raining.*
 - **Probable:** *It is **probably** raining.*

- **Hypothetical:**

- This type of uncertainty includes four sub-classes:
 - **Doxastic:** I **believe** Tom can win the game.
 - **Investigation:** I **examined** the result and found
 - **Condition:** **If** tom can win, I will buy you lunch.
 - **Dynamic:** I **hope** tom can win.

*Ferenc Kiefer. 2005. *Lehetoseg es szuksegszeruseg* [Possibility and necessity]. Tinta Kiado, Budapest.



Preliminary experiment

- 827 tweets annotation
 - Traditional scheme: 65 uncertain
 - Manually: 246 uncertain
 - More than **70%** uncertain tweet are missing.
- *Different uncertainty expression on social media.*



Uncertainty in social media

- Three observations

- No tweet under category of **investigation**.
 - @dobibid I have tested the link, it is fake!
- Express uncertainty by **question**.
 - @ITVCentral *Can you confirm that Birmingham children's hospital has/hasn't been attacked by rioters?*
- Express uncertainty by quoting **external** information.
 - *Friend who works at the children's hospital in Birmingham says the riot police are protecting it.*



Classification for social media

Category	Subtype	Cue	Example
Epistemic	Possible	may	It may be raining.
	Probable	likely	It is probably raining.
Hypothetical	Condition	if	If it rains, we'll stay in.
	Doxastic	believe	He believes that the Earth is flat.
	Dynamic	hope	fake picture of the london eye on fire... i hope
	External	someone said	Someone said that London zoo was attacked.
	Question	seriously?	Birmingham riots are moving to the children hospital?! seriously?

- Based on proposed scheme is based on Kiefer's work (2005) which was previously extended to normalize uncertainty corpora in different genres by Szarvas et al. (2012).
- Ferenc Kiefer. 2005. *Lehetoseg es szuksegszeruseg*[Possibility and necessity]. Tinta Kiado, Budapest.
- György Szarvas, Veronika Vincze, Richárd Farkas, György Móra, and Iryna Gurevych. 2012. Crossgenre and cross-domain detection of semantic uncertainty. *Computational Linguistics*, 38(2):335–367.



Annotation

- London Riots dataset
 - August 6-13 2011
 - 4,743 unique tweet related to seven riots events*.
- Annotation scheme
 - Two trained annotators.
 - Binary judgment in terms of author's intended meaning.
 - Sub-class label for tweets with uncertainty label.
 - A third annotator for final decision.
 - Cue-phrase identification to form a uncertainty cue-phrase list.

Identified by UK newspaper “*The Guardian*”



Annotation

- Tweet #: 4743
- Uncertainty#: 926 (19.52%)
- Kappa agreement:
 - 0.9073 for binary classification
 - 0.8271 for fine-grained annotation

Epistemic	Possible#	16
	Probable#	129
Hypothetical	Condition#	71
	Doxastic#	48
	Dynamic#	21
	External#	208
	Question#	488



Experiment setup

- Task
 - Uncertainty tweet identification
- Approaches
 - Cue-phrase matching (CP)
 - Supervised machine learning (SVM_{***})
 - N-grams (unigram + bigram + trigram)
 - Content-based feature
 - Twitter-specific feature
 - User-based feature
- Evaluation
 - 5-fold validation
 - Precision, recall, F-1



Experiment

Category	Name	Description
Content-based	Length	Length of the tweet
	Cue_Phrase	Whether the tweet contains a uncertainty cue
	OOV_Ratio	Ratio of words out of vocabulary
Twitter-specific	URL	Whether the tweet contains a URL
	URL_Count	Frequency of URLs in corpus
	Retweet_Count	How many times has this tweet been retweeted
	Hashtag	Whether the tweet contains a hashtag
	Hashtag_Count	Number of hashtag in tweets
	Reply	Is the current tweet a reply tweet
	Retweet	Is the current tweet a retweet tweet
User-based	Follower_Count	Number of follower the user owns
	List_Count	Number of list the users owns
	Friend_Count	Number of friends the user owns
	Favorites_Count	Number of favorites the user owns
	Tweet_Count	Number of tweets the user published
	Verified	Whether the user is verified



Experiment

Approach	Precision	Recall	F-1
CP	0.3732	0.9589	0.5373
SVM _{n-gram}	0.7278	0.8259	0.7737 (+43.9%*)
SVM _{n-gram+C}	0.8010	0.8260	0.8133
SVM _{n-gram+U}	0.7708	0.8271	0.7979
SVM _{n-gram+T}	0.7578	0.8266	0.7907
SVM _{n-gram+ALL}	0.8162	0.8269	0.8215

- C: content based features.
- U: user based features.
- T: twitter specific features.
- ALL: the combination of C, U and T.

***compare to CP**



Experiment

- Performance of content-based features

Approach	Precision	Recall	F-1
SVM _{n-gram+Cue-Phrase}	0.7989	0.8266	0.8125
SVM _{n-gram+Length}	0.7372	0.8216	0.7715
SVM _{n-gram+OOV_Ratio}	0.7414	0.8233	0.7802

- Presence of uncertain cue-phrase is most indicative.



Experiment

- Classification errors of SVM_{n-gram+ALL}

Type	Poss.	Prob.	D.&D.	Cond.	Que.	Ext.
Total#	16	129	69	71	488	208
Error#	11	20	18	11	84	40
Error%	0.69	0.16	0.26	0.15	0.17	0.23

- Combine dynamic and doxastic for error analysis.
- Perform worst on two categories with least samples.



Conclusion

- Propose a variant of classification scheme for uncertainty identification in social media.
- Perform uncertainty identification experiments and explore effectiveness of different type of features.
- In future, we will explore to use uncertainty identification for social media applications



Questions or Suggestions?



Zhongyu Wei (魏忠鈺)

[http://www.se.cuhk.edu.hk/~zywei/
zywei@se.cuhk.edu.hk](http://www.se.cuhk.edu.hk/~zywei/zywei@se.cuhk.edu.hk)

Kam-Fai Wong(黃錦輝)

[http://www.cintec.cuhk.edu.hk/kfwong/
kfwong@se.cuhk.edu.hk](http://www.cintec.cuhk.edu.hk/kfwong/kfwong@se.cuhk.edu.hk)

