

Re-tweeting from a Linguistic Perspective

Aobo Wang, Tao Chen and Min-Yen Kan

7/06/2012

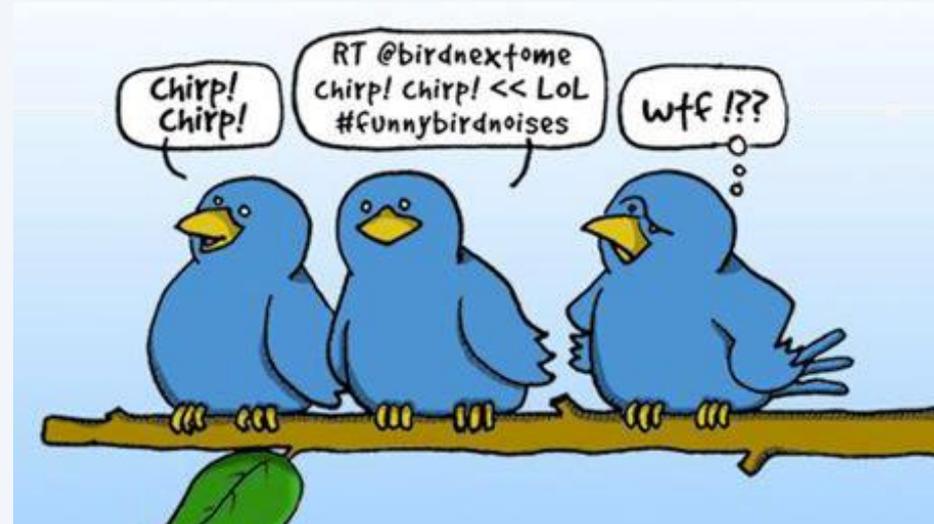
Introduction



NUS Libraries @NUSLibraries

Prosperous state, prosperous old? :growing social stratification among elderly Singaporeans - free "Ebook" from ARI j.mp/Mijl5c

Expand Reply Retweet Favorite



Q: What makes a tweet worth sharing?

- from a linguistic perspective

Introduction

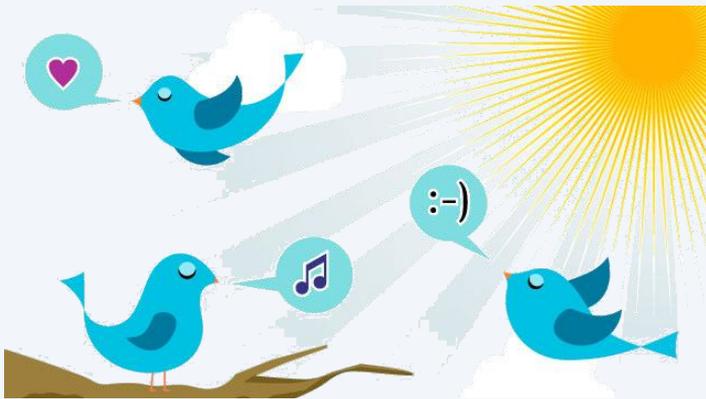
- **Something we know**



- **Social network effects** exert marked influence on retweeting
(Wu et al., 2011; Recuero et al., 2011)

Motivation

- Something we want to know



Q: Are there specific **linguistic signals** that mark a tweet as valuable and worthy of sharing?

Tasks

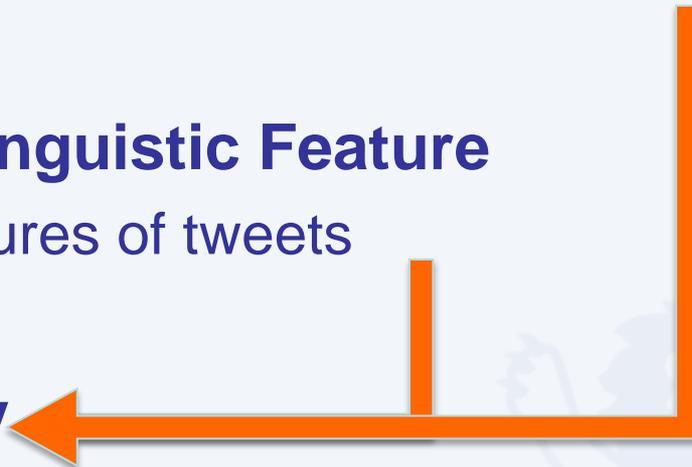
1. Linguistically Motivated Tweet Classification

- The specific function of the individual tweet

1. Analysis of Linguistic Feature

- Linguistic features of tweets

1. Retweetability



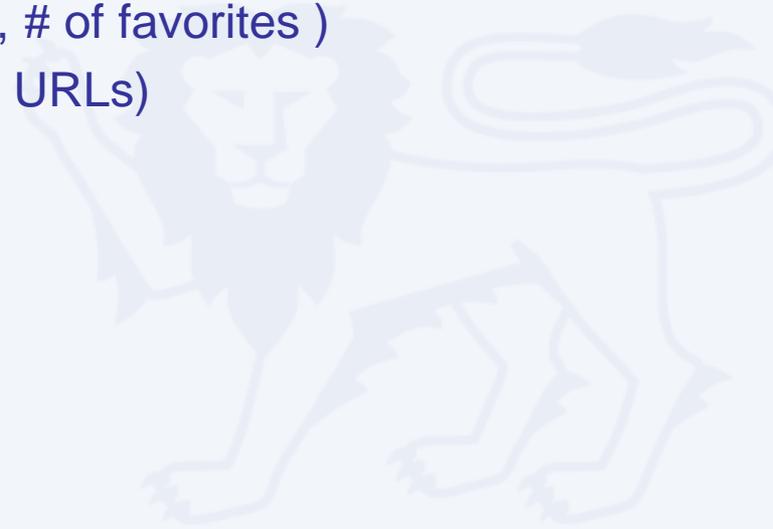
Literature Review

- **Automatic Classification**

- Sriram et al. 2010

- 5 genres classification scheme
 - Supervised method using Naïve Bayes Classifier
 - 5407 manually labeled tweets
 - Domain-specific features from
 - **author's profile** (e.g., # of followers, # of favorites)
 - **lexicon of tweets** (e.g., #hashtags, URLs)
 - **metadata** (time phrases).

Opinions
Events
News
Deals
Private Messages



Literature Review

- **Automatic Classification**

- Ramage et al. 2010

- Semi-supervised method, indirect tweet level classification

1. Unsupervised labelling tweets with topic label

- Get topic labels with LDA as **Topic Set A**

- Treat **Hashtag**, **Emoticons**, and **Social Signal (@user)** as **Topic Set B**

2. Manually classify the **Set A+B** into 4 genres.

3. Train Labeled LDA classification model with the **Set A+B** topic labels

- **We know little about the linguistic features of tweets.**

- **Classify tweets based on the functions of tweets using linguistic features.**

Hypothesis

- Tweets with particular **function** will be used when users have corresponding **motivations** of tweeting.
- People's **motivations** in posting tweets determine their **writing styles**.
- Such **styles** can be characterized by the **content** and **linguistic features** of tweets.
 - “I am presenting in Salon now.”



Data Set Collection

- **More than 9 million tweets crawled by Twitter Stream API**
- **Pre-processing**
 - Exclude tweets with URLs from our current study
 - Break the hashtags into separate words
(e.g., #growingup → growing up)
 - Tokenizing on emoticons, usernames (@user) and “RT if”-like (“retweet if”) syntax patterns.



Data Annotation

- Classification scheme and Example tweets

Level-1	Level-2	Motivation	Example	Corpus count (%)
Opinion	Abstract	Present opinions towards abstract objects	<i>God will lead us all to the right person for our lives. Have patience and trust him.</i>	291 (33.8%)
	Concrete	Present opinions towards concrete objects	<i>i feel so bad for nolan. Cause that poor kid gets blamed for everything, and he's never even there.</i>	99 (11.5%)
	Joke	Tell jokes for fun	<i>Hi. I'm a teenager & I speak 3 languages: English, Sarcasm, & Swearing (; #TeenThings</i>	86 (10.0%)
Update	Myself	Update my current status	<i>first taping day for #growingup tomorrow! So excited. :)</i>	168 (19.6%)
	Someone	Update others' current status	<i>My little sister still sleep ...</i>	66 (7.7%)
Interaction		Seek interactions with others.	<i>#Retweet If you're #TeamFollowBack</i>	81 (9.4%)
Fact		Transfer information	<i>Learnt yesterday: Roman Empire spent 75% of GDP on infrastructure. Roads, aqueducts, etc.</i>	23 (2.7%)
Deals		Make deal	<i>Everybody hurry! Get to Subway before they stop serving LIMITED TIME ONLY item 'avocados'.</i>	29 (3.4%)
Others		Other motivations.	<i>Ctfu Lmfao At Kevin Hart ;)</i>	17 (2.0%)

- Collect Labels through Amazon's Mechanical Turk

- 860 tweets in total
- Fleiss' kappa : Level-1=0.79; Level-2=0.43



Method

- **Labeled LDA Classification**
 - Tweet level classification on **Level 1**
 - 5-fold validation
 - Feature selection
 - **Content**
 - **Discourse relations**
 - **Hashtags**
 - **Interaction Lexical patterns**
 - **Named Entities**
 - **Tense**
- **Incremental training**

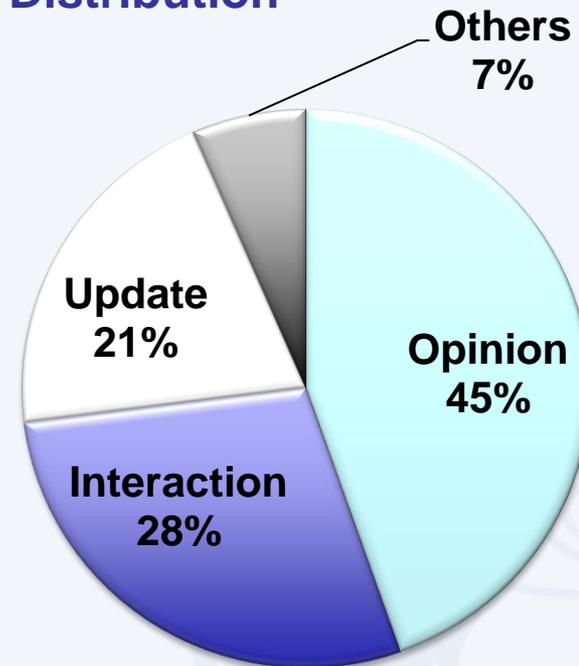


Classification Result

- Weighted average F-1 Score

	Level-1	Level-2
C (baseline)	.625	.413
CD	.637	.432
CH	.629	.415
CI	.642	.422
CN	.611	.409
CT	.635	.427
CDHIT	.670	.451

- Distribution



Tasks

1. Linguistically Motivated Tweet Classification

- The specific function of the individual tweet
- More than 9 million tweets

1. Analysis of Linguistic Features

- Linguistic features of tweets
- More than 1.5 million retweets

1. Retweetability



Emoticons and Sentiment

- :) → positive :(→ negative
 - Read et al. 2005, Go et al. 2009, Alexander et al. 2010

➤ **Q: Do emoticons actually indicate sentiment of message?**

- Randomly select 200 posts with smilies and 200 posts with frownies
- Label their sentiment manually
- Evaluate Go et al. (2009)’s API on our annotated corpus

	Positive	Neutral	Negative
Retweets with :)	55 (27.5%)	140 (70%)	5 (2.5%)
Retweets with :(9 (4.5%)	118 (59%)	73(36.5%)
Predicted Positive	43	30	0
Predicted Neutral	11	206	12
Predicted Negative	7	29	62

Majority is neutral tweets

Mistake neutral posts for emotional ones

➤ **Use emoticons carefully to detect sentiment**

Named Entities

- **Q: What types of NEs do people mention in their tweets?**
 - Extract NEs by UW Twitter NLP Tools (Ritter et al., 2011)
 - Select the top 100 correctly recognized NEs
 - Manually categorize NEs against their 10 classes scheme (defined by Ritter et al. 2011)

Class	Opinion	Update	Interaction
<i>PERSON</i>	41.2%	44.7%	38.8%
<i>GEO-LOC</i>	7.8%	28.9%	25.4%
<i>COMPANY</i>	15.7%	6.6%	10.4%
<i>PRODUCT</i>	5.9%	5.3%	6.0%
<i>SPORTS-TEAM</i>	2.0%	5.3%	1.5%
<i>MOVIE</i>	7.8%	5.3%	7.5%
<i>TV-SHOW</i>	3.9%	0.0%	3.0%
<i>OTHER</i>	15.7%	3.9%	7.5%

- **Person Names are dominating.**
- **Geographical locations are less often mentioned in Opinion**

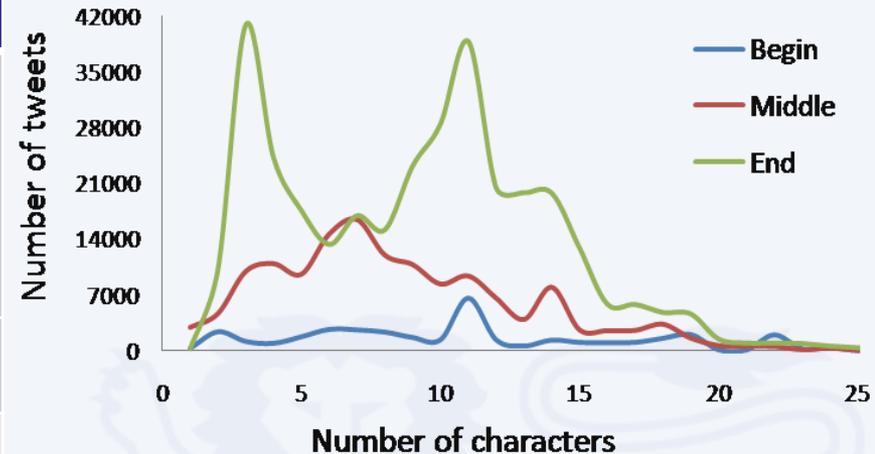


Hashtags

Q: Any positional preference for embedding hashtags?

Position	Example Tweets	%
End	Success is nothing without someone you love to share it with. #TLT	69.1
	Goodmorning Tweethearts....wishing u all blessed and productive day! #ToyaTuesday	
Middle	I just saw the #Dodgers listed on Craig's List.	20.7
Beginning	#ihateit when random people poke you on facebook	8.9

Q: Any patterns to how people form hashtags?

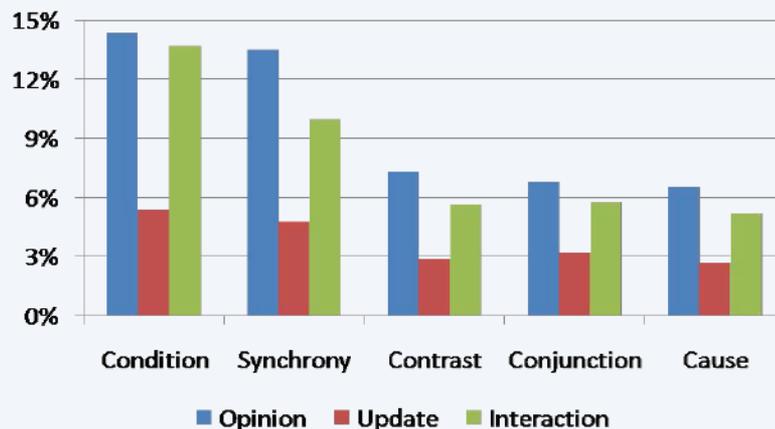


- **Enders:** peak at around 3 or 11 → Twitter slang, time and location
- **Middlers:** peak at around 7 → Single keyword
- **Beginners:** peak at around 11 → subject+verb+object

Discourse Relation and Sentence Similarity

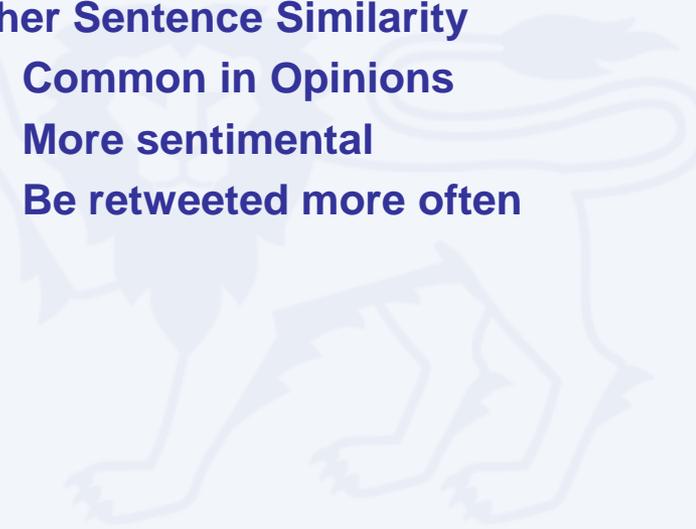
- **Discourse Relation**
 - End-to-end discourse parser by Lin et al. (2010)
 - PDTB-styled discourse relations (Prasad et al. 2008)

➤ **Five most frequent relations**



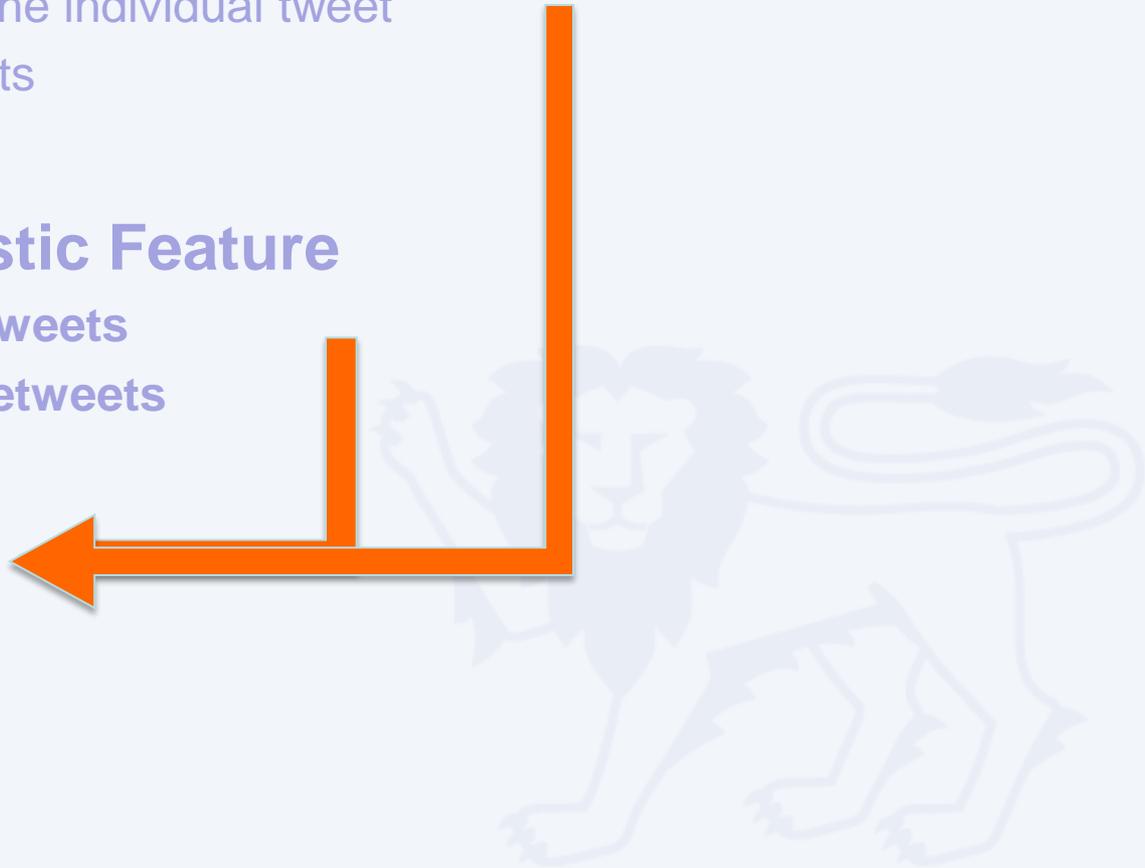
- **Sentence Similarity**
 - Example: *“On Twitter people follow those they wish they knew. On Facebook people follow those they used to know.”*
 - Computed by Syntactic Tree Matching model (Wang et al. 2009)

- **Higher Sentence Similarity**
 - **Common in Opinions**
 - **More sentimental**
 - **Be retweeted more often**



Tasks

- **Linguistically Motivated Tweet Classification**
 - The specific function of the individual tweet
 - More than 9 million tweets
- **Analysis of Linguistic Feature**
 - Linguistic features of tweets
 - More than 1.5 million retweets
- **Retweetability**



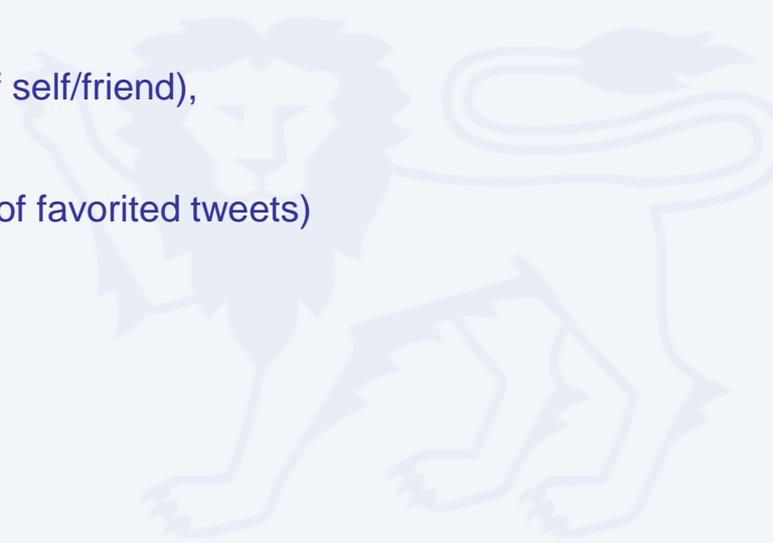
Literature Review

- **Previous work**

- Retweet rate prediction using GLM; Suh et al., (2010)
- Retweet probability prediction using CRF; Peng et. Al (2011)
- Retweet volume prediction using Logistic Regression; Hong et al.(2011)

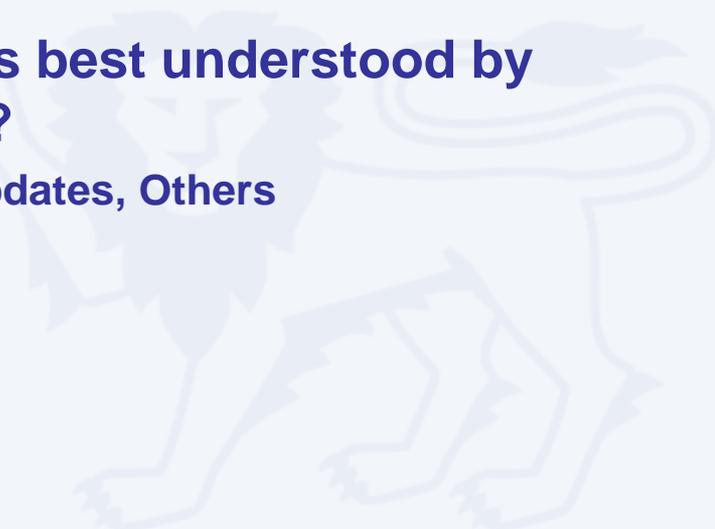
- **Previously Examined Feature Sets**

- Author's profile
 - (e.g., # of followers/followees/friends; activity of self/friend),
- Tweet metadata
 - (e.g., time interval,# of previously retweeted, # of favorited tweets)
- Twitter-specific features
 - (URL , Hashtags, @user)



What does the tweet itself contribute to its retweetability?

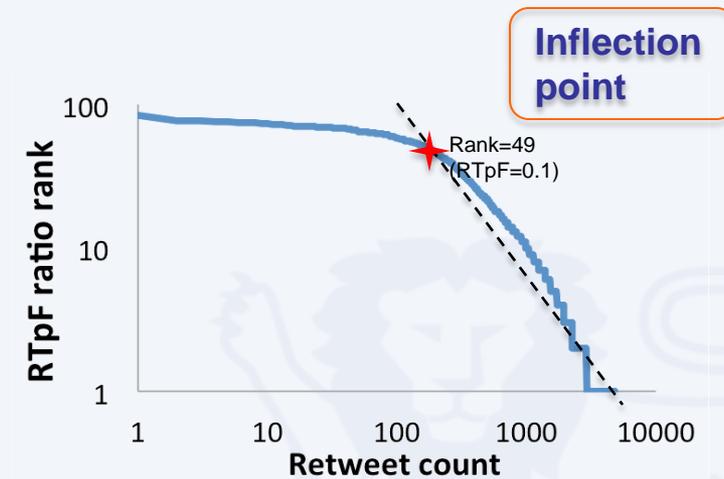
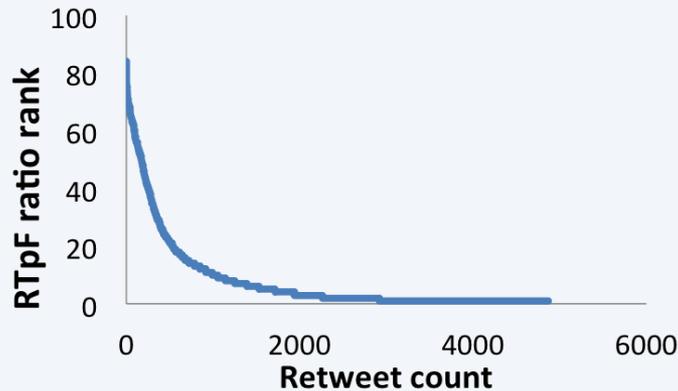
- Surface level features
 - Presence of hashtags, @user, quotation, 3 hashtag positions
 - Tweet length, hashtag counts
- Linguistic features
 - Presence of 16 types of discourse relations; 10 NE types; Verb tenses; 3 sentiment polarity strengths
 - Sentence similarity value
- **Whether a tweet is shared with others is best understood by modeling each function independently?**
 - **Level-1** functions: **Opinion, Interaction, Updates, Others**
- **Tweet content is not factored**



Experiment

- Task Definition**

- $RTpF = \# \text{ of Retweets} / \text{Followers count}$
- Given the content of a tweet, perform a multi-class classification that predicts its range of **RTpF** ratio.



- Non-retweets (“N”, $RTpF = 0$),
- Low (“L”, $RTpF < 0.1$),
- High (“H”, $RTpF > 0.1$)

Experiment

- **Data Set**
 - Selected from 9 million dataset
 - Balanced data size of three RTpF classes.
- **Method**
 - Logistic Regression model in Weka3
 - 10-fold cross validation



Result

- Individual regression models
- Aggregate models for all three classes

Class	F_1
<i>Opinion</i>	0.57
<i>Update</i>	0.54
<i>Interaction</i>	0.53
All w/o L-1 class	0.42
All w/ L-1 class	0.52

Independent models perform better than combined model

The usage of Level-1 feature improves performance

Observation and Remarks

- Opinion

Salient Features	Weight	Example Tweets	RTpF
Sentence Similarity	10.34	<i>“twitter is where people vent to vent, facebook is where people vent to get attentionn”</i>	0.84
Conjunction	-21.09	<i>“#Cancer #Scorpio and #Pisces will become quiet and withdrawn when things get tough and they need to think.”</i>	0.10
Quotation	-19.2	<i>“If you obey all the rules, you miss all the fun - Katharine Hepburn”</i>	0.22

- Beautiful sentence structure
- Avoid complex conjoined components
- Make your words originally



Observation and Remarks

- **Update**

Salient Features	Weight	Example Tweets	RTpF
Past	-5.2	<i>“I fell for your personality, and your looks were just a bonus”</i>	0.08
Present	1.3	<i>“Lying in bed, wondering if its worth it to get up”</i>	0.17

- Shows the least bias towards any particular feature
- Prefers present tenses to past tense



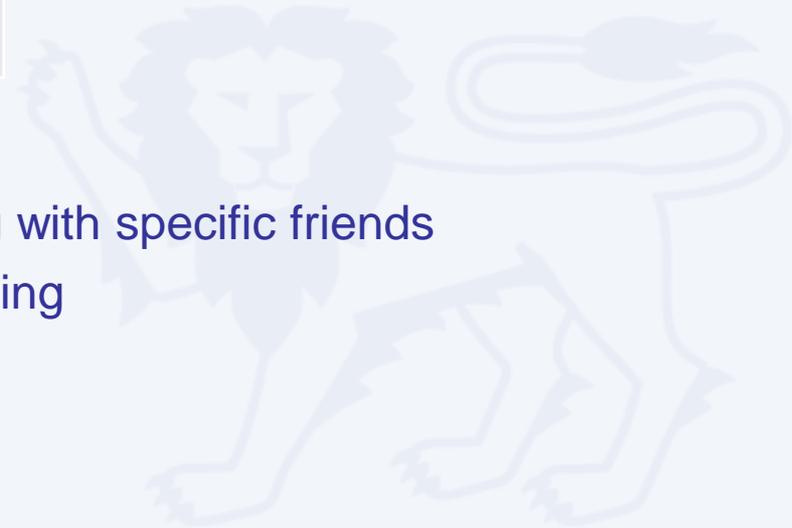
Observation and Remarks

- **Interaction**

–“--->**R E T W E E T**<--- *If you want more followers #TeamFollowBack | #TFB | #InstantFollowBack | #500ADay | #MustFollow @iTweetHeavyTGOD*”

Salient Features	Weight
Sentence Similarity	-55.33
Hashtag Count	5.34

- Keep direct and simple while interacting with specific friends
- In the form of question answering or voting



Observation and Remarks

- **Globally**

Class	Salient Feature	Weight
All w/o -1 class	Hashtag Count	22.03
All w/ L-1 class	Sentence Similarity	9.8

- Hashtags are positive triggers
- L-1 Class features are important



Conclusion

- Understanding and classify the function of the tweet is interesting in its own right.
- It is also useful in predicting the retweetability.
- **Release**
 - A corpus of 860 annotated tweets
 - Functional classifier
 - Online demo
 - <http://wing.comp.nus.edu.sg/tweets/>
- Tweets containing URLs and the features from social network perspective will be taken into consideration in future work.



Re-tweeting from a Linguistic Perspective

Thank you very much!

Aobo Wang, Tao Chen and Min-Yen Kan

07/06/2012

