

Good Friends, Bad News

*How emotions influence
virality in Twitter*

Lars Kai Hansen

Dept of Signal Theory and
Communications,

Universidad Carlos III de Madrid,
28911 Leganes, Spain

www.imm.dtu.dk/~lkh



DTU Informatics

Department of Informatics and Mathematical Modeling

Thanks to my collaborators in the
Responsible business in the blogosphere project

Good Friends, Bad News - Affect and Virality in Twitter

Lars Kai Hansen, Adam Arvidsson, Finn Aarup Nielsen,
Elanor Colleoni, and Michael Etter

DTU Informatics, Technical University of Denmark,
DK-2800 Lyngby, Denmark, Email: lkh@imm.dtu.dk fn@imm.dtu.dk

University of Milan, via Conservatorio Milan, Italy,

Email: adam.arvidsson@unimi.it

Copenhagen Business School, DK-2000 Frederiksberg, Denmark

Email: elc.ikl@cbs.dk me.ikl@cbs.dk



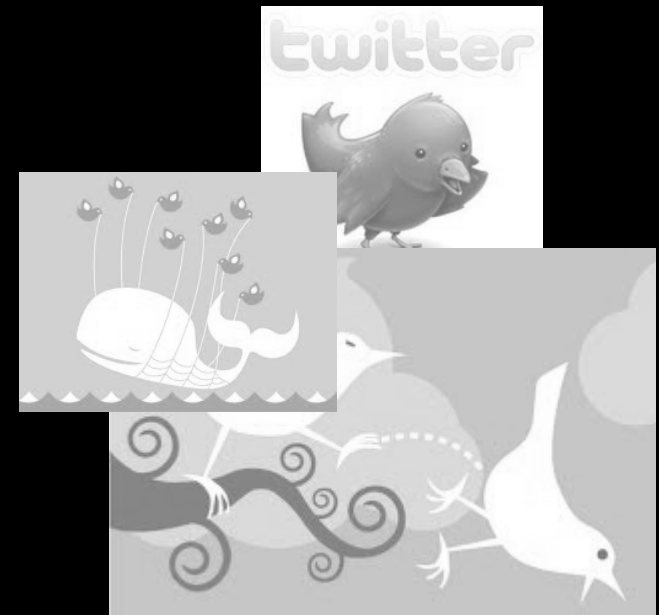
In press: SocialComNet 2011 (Athens, June 28, 2011)



Outline



- Social media outlets are instrumental for the new, data driven science of human behavior
 - Computers can read (starting to understand)
 - Quantifying subjectivity: Affective computing
 - Sentiment analysis (in text)
- The sentiment/virality paradox
 - Social interaction is positive
 - Twitter as a laboratory
 - Analyzing virality in Twitter
 - Paradox resolved
- Conclusions



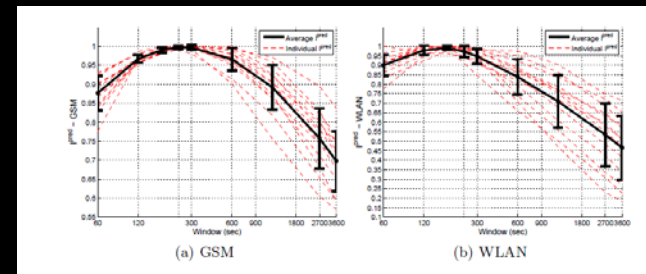


Human predictability

C. Song, Z. Qu, N. Blumm, A.-L. Barabási
 Limits of Predictability in Human Mobility
 Science 327, 1018-1021 (2010).

A new picture of human predictability is emerging

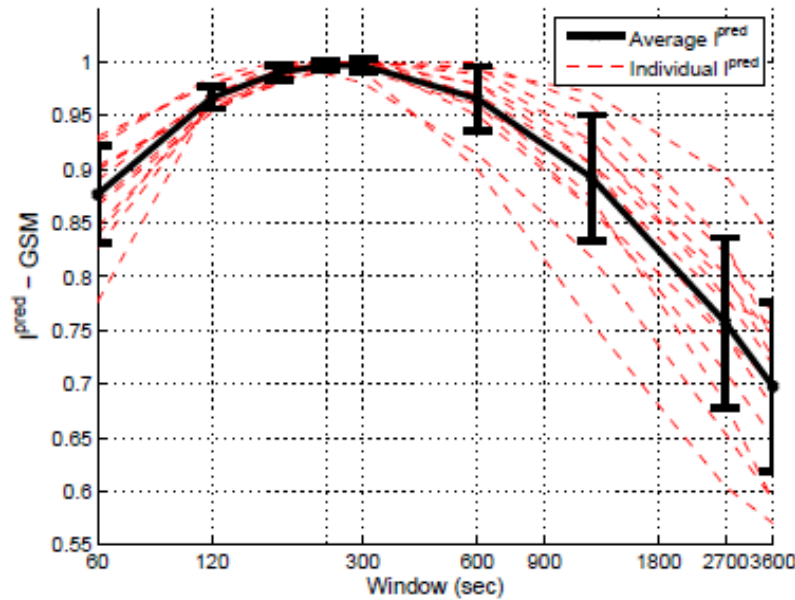
- Barabasi group's study of mobility (70-93% /hour)
- Individuals are predictable (intersubject variability is huge ... long tail ... "power law")
- DTU mobile phone context tracking (80-90% /4 min)



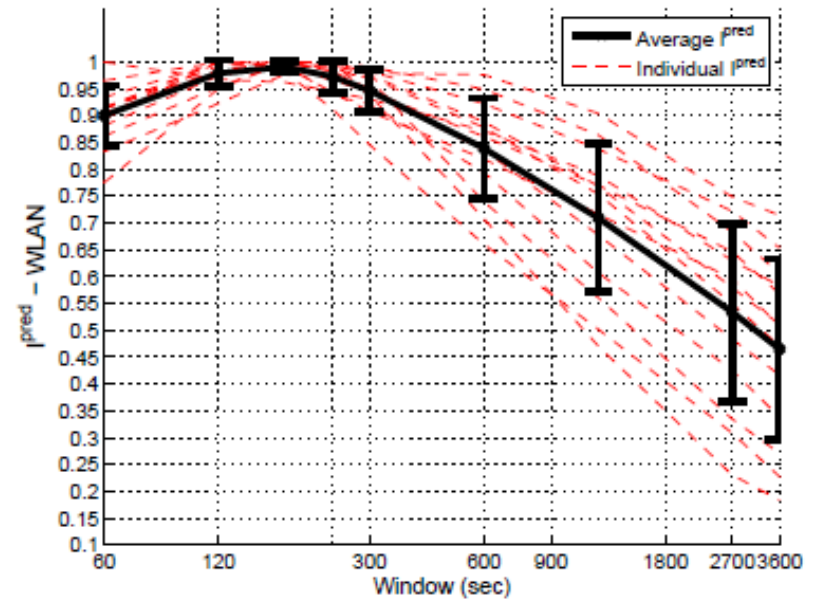
Why is this important?

- **Curiosity:** Understanding the human brain
- **Engineering:** New services/products are based on human predictability (the Google paradigm)

Human predictability

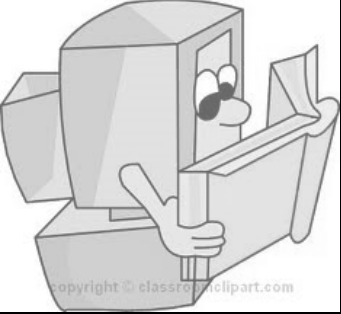


(a) GSM



(b) WLAN

B.S. Jensen, J. Larsen, K. Jensen, J.E: Larsen, L.K. Hansen
Preditability of mobile phone associations. In Proc MUSE (2010)



Computers can read (& are starting to understand)

Main approaches to text analysis in computers

- Computational linguistics (NLP)
 - Rule based, tight statistical models
 - Precise / high specificity
 - Challenged by informal/innovative text
- Statistical learning
 - Flexible / unsupervised
 - High sensitivity -informal OK
 - Challenged in the quantitative/opinion/precision



Trick: Vector space representation

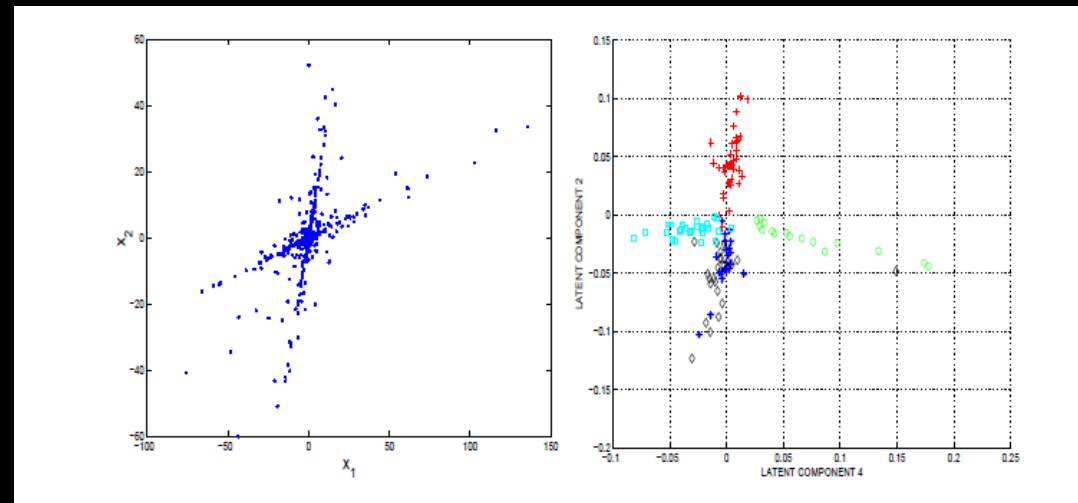
- Abstract representation - can be used for all digital media
- Document is represented as a point in a high-dimensional "feature space" - document similarity \sim spatial proximity
- General features or "events"
- Social network: Social event involving a set of nodes
- Text: Bag of words (Term/keyword histograms),
- Image: Color histogram, texture measures, "bag of features"
- Video: Object coordinates (tracking), active appearance models
- Sound: Spectrograms, cepstral coefficients, gamma tone filters

Document features are correlated, the pattern of correlation reflects "associations". Associations are context specific

Contexts can be identified unsupervised fashion by their feature associations = Latent semantics

"Bag of Words"

Terms	Documents									
	c1	c2	c3	c4	c5	m1	m2	m3	m4	
computer	1	1	0	0	0	0	0	0	0	
EPS	0	0	1	1	0	0	0	0	0	
human	1	0	0	1	0	0	0	0	0	
interface	1	0	1	0	0	0	0	0	0	
response	0	1	0	0	1	0	0	0	0	
system	0	1	1	2	0	0	0	0	0	
time	0	1	0	0	1	0	0	0	0	
user	0	1	1	0	1	0	0	0	0	
graph	0	0	0	0	0	0	1	1	1	
minors	0	0	0	0	0	0	0	1	1	
survey	0	1	0	0	0	0	0	0	1	
tree	0	0	0	0	0	1	1	1	0	

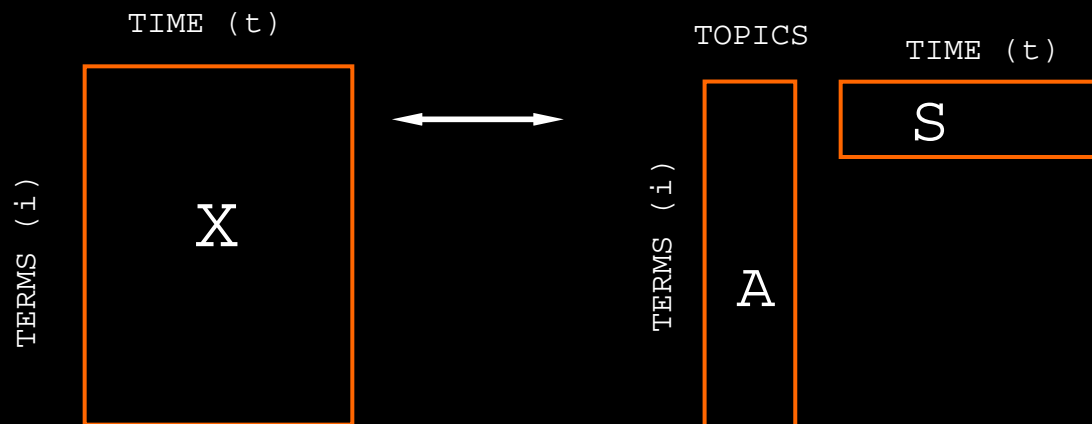


Very efficient for detection of context

Often leads to very high-dimensional learning problems

Factor models

- Represent a datamatrix by a low-dimensional approximation

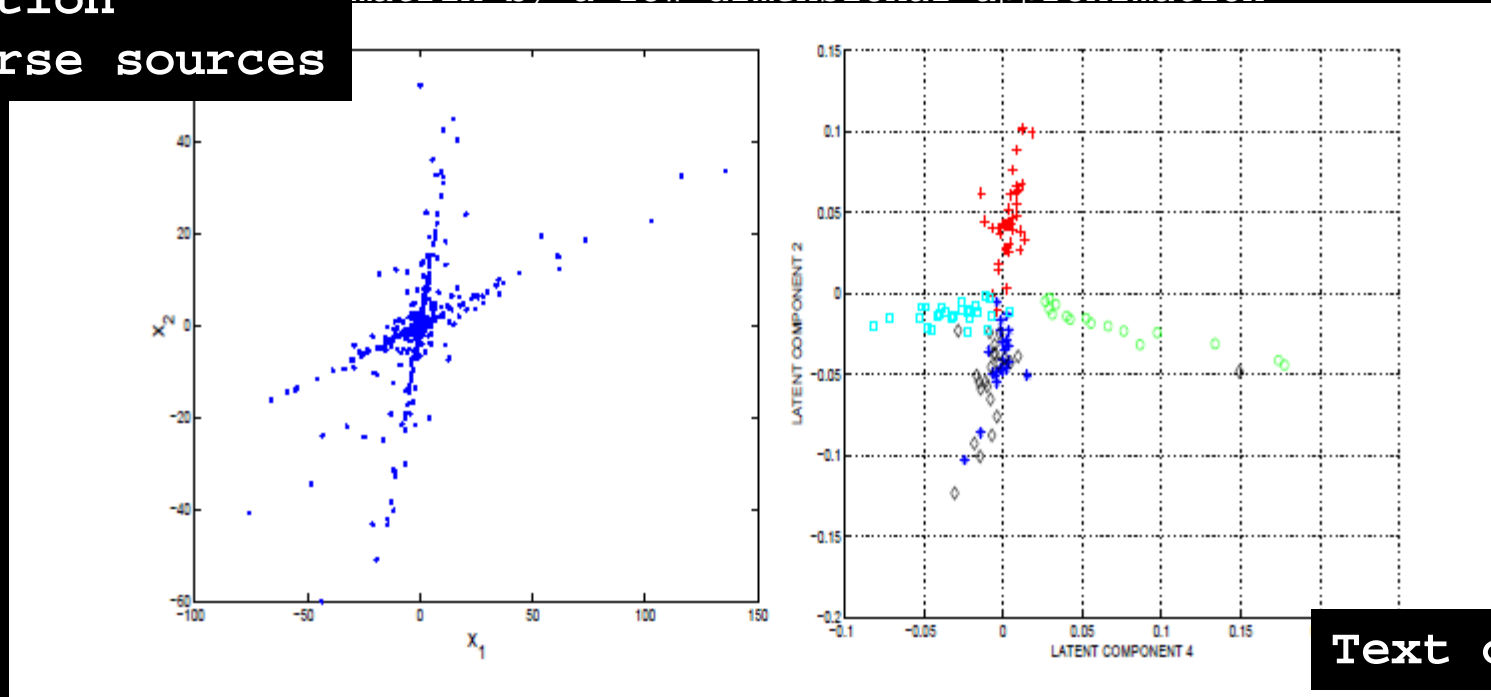


$$X(i, t) \approx \sum_{k=1}^K A(i, k) S(k, t)$$

Factor models

Simulation
w/ sparse sources

matrix by a low-dimensional approximation



$$X(i, t) \approx \sum_{k=1}^K A(i, k) S(k, t)$$

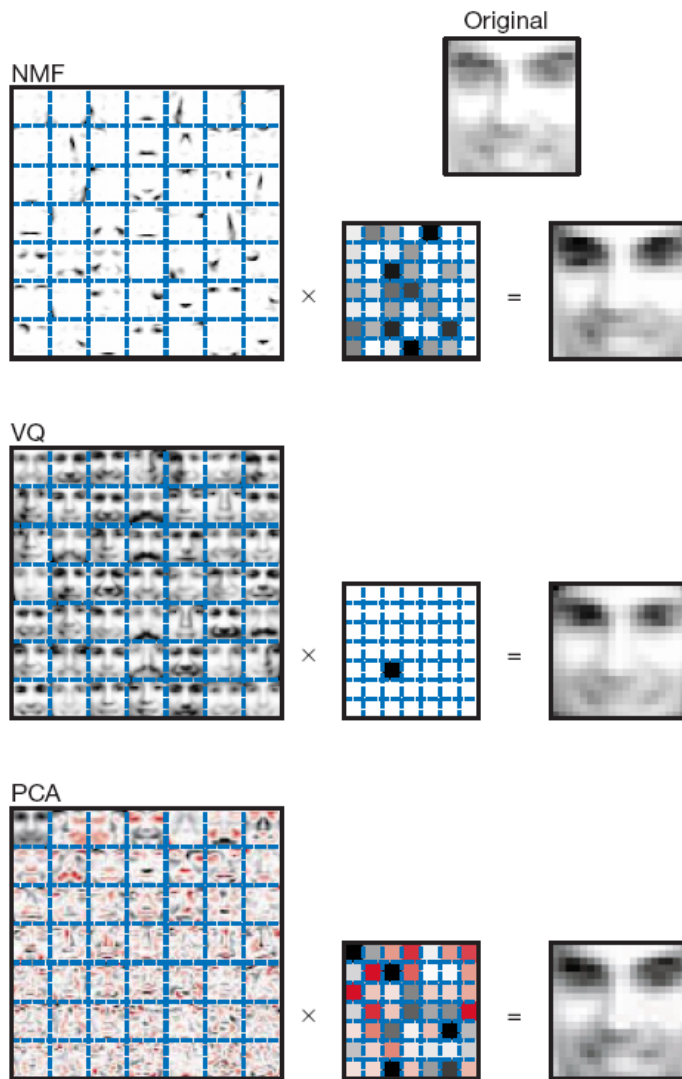


Figure 1 Non-negative matrix factorization (NMF) learns a parts-based representation of faces, whereas vector quantization (VQ) and principal components analysis (PCA) learn holistic representations. The three learning methods were applied to a database of $m = 2,429$ facial images, each consisting of $n = 19 \times 19$ pixels, and constituting an $n \times m$ matrix V . All three find approximate factorizations of the form $V \approx WH$, but with three different types of constraints on W and H , as described more fully in the main text and methods. As shown in the 7×7 montages, each method has learned a set of $r = 49$ basis images. Positive values are illustrated with black pixels and negative values with red pixels. A particular instance of a face, shown at top right, is approximately represented by a linear superposition of basis images. The coefficients of the linear superposition are shown next to each montage, in a 7×7 grid, and the resulting superpositions are shown on the other side of the equality sign. Unlike VQ and PCA, NMF learns to represent faces with a set of basis images resembling parts of faces.

Learning the parts of objects by non-negative matrix factorization

Daniel D. Lee* & H. Sebastian Seung*†

* Bell Laboratories, Lucent Technologies, Murray Hill, New Jersey 07974, USA

† Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA

NATURE | VOL 401 | 21 OCTOBER 1999 | www.nature.com

CASTSEARCH - CONTEXT BASED SPEECH DOCUMENT RETRIEVAL

Lasse Lohilahti Mølgaard, Kasper Winther Jørgensen, and Lars Kai Hansen

Informatics and Mathematical Modelling
Technical University of Denmark Richard Petersens Plads
Building 321, DK-2800 Kongens Lyngby, Denmark

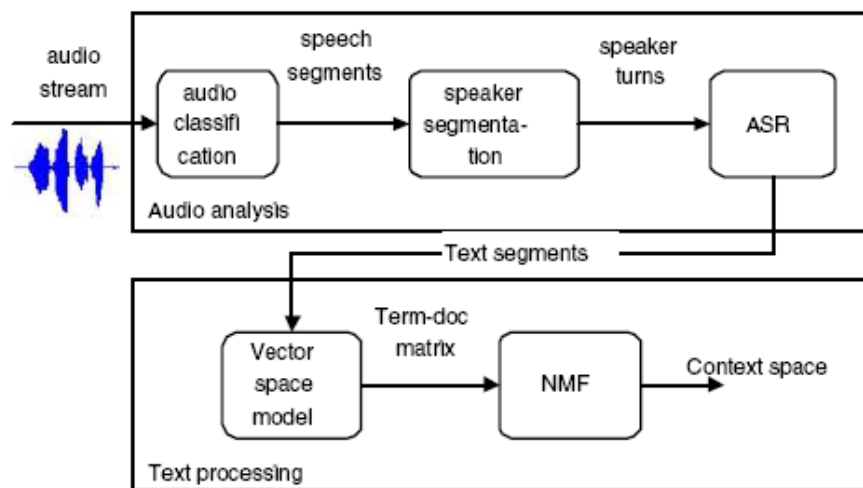
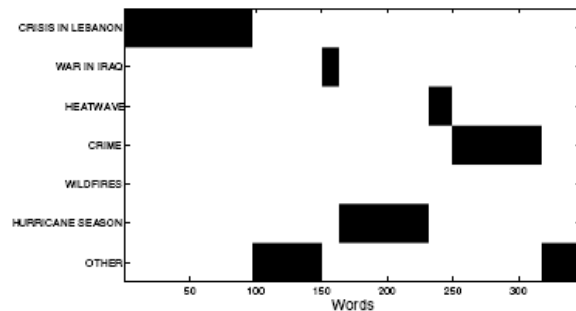


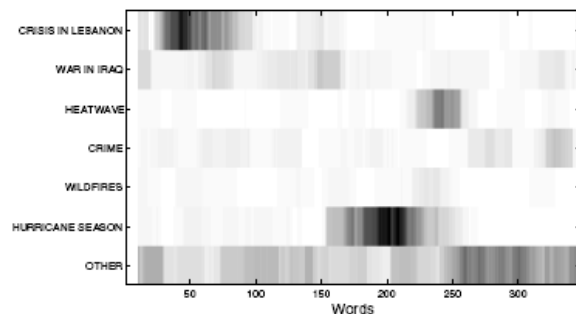
Fig. 1. The system setup. The audio stream is first processed using audio segmentation. Segments are then using an automatic speech recognition (ASR) system to produce text segments. The text is then processed using a vector representation of text and apply non-negative matrix factorization (NMF) to find a topic space.

DTU Informatics

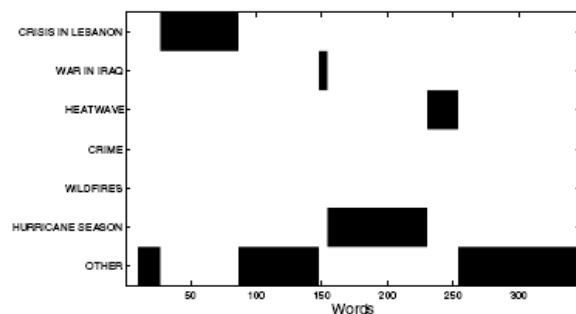
Department of Informatics and Mathematical Modeling



(a) Manual segmentation.

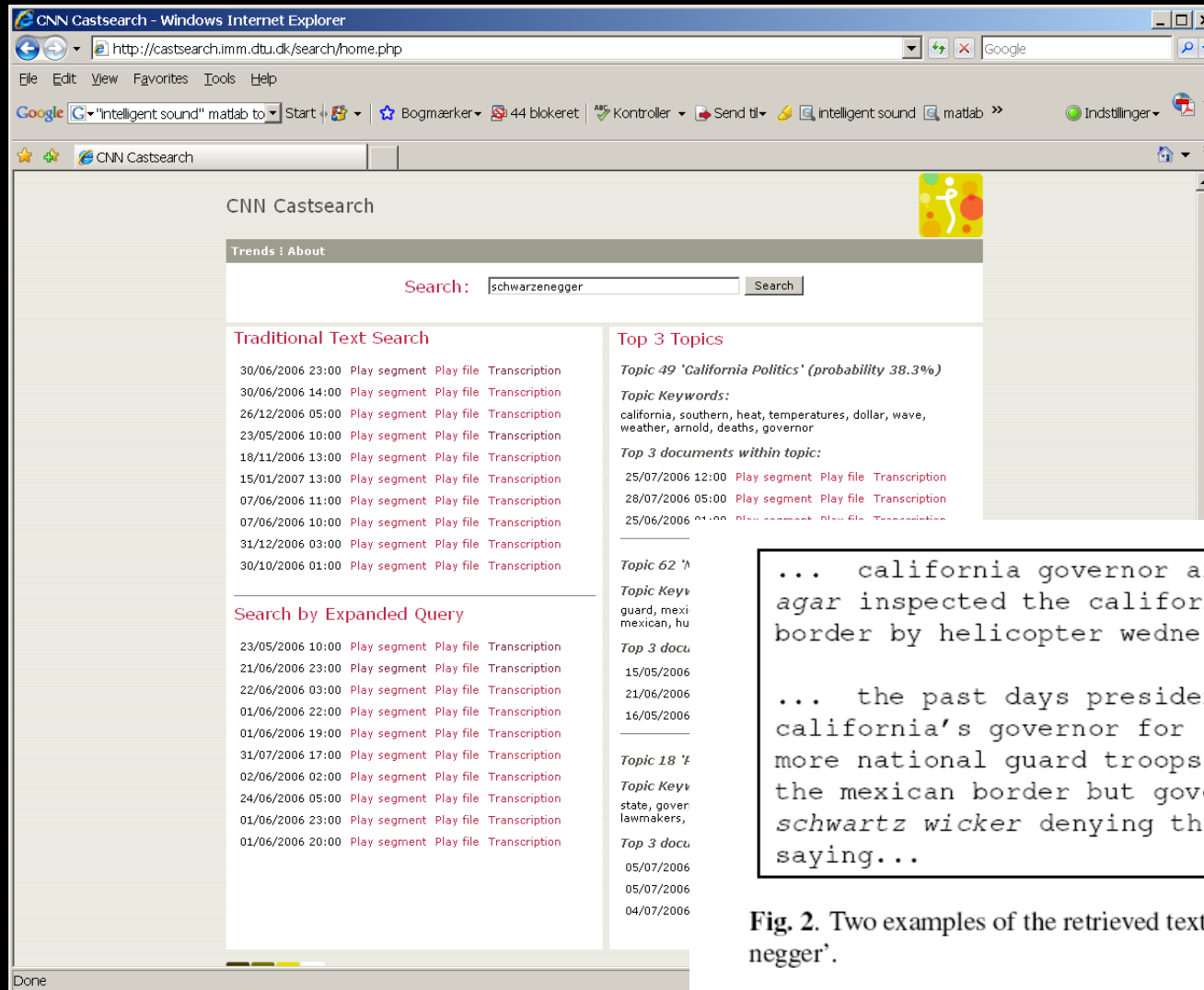


(b) $p(k|d^*)$ for each context. Black means high probability.



(c) The segmentation based on $p(k|d^*)$.

Fig. 3. Figure 3(a) shows the manual segmentation of the news show into 7 classes. Figure 3(b) shows the distribution $p(k|d^*)$ used to do the actual segmentation shown in figure 3(c). The NMF-segmentation is in general consistent with the manual segmentation. Though, the segment that is manually segmented as 'crime' is labeled 'other' by the NMF-segmentation



CNN Castsearch

Trends : About

Search: Search

Traditional Text Search

30/06/2006 23:00	Play segment	Play file	Transcription
30/06/2006 14:00	Play segment	Play file	Transcription
26/12/2006 05:00	Play segment	Play file	Transcription
23/05/2006 10:00	Play segment	Play file	Transcription
18/11/2006 13:00	Play segment	Play file	Transcription
15/01/2007 13:00	Play segment	Play file	Transcription
07/06/2006 11:00	Play segment	Play file	Transcription
07/06/2006 10:00	Play segment	Play file	Transcription
31/12/2006 03:00	Play segment	Play file	Transcription
30/10/2006 01:00	Play segment	Play file	Transcription

Search by Expanded Query

23/05/2006 10:00	Play segment	Play file	Transcription
21/06/2006 23:00	Play segment	Play file	Transcription
22/06/2006 03:00	Play segment	Play file	Transcription
01/06/2006 22:00	Play segment	Play file	Transcription
01/06/2006 19:00	Play segment	Play file	Transcription
31/07/2006 17:00	Play segment	Play file	Transcription
02/06/2006 02:00	Play segment	Play file	Transcription
24/06/2006 05:00	Play segment	Play file	Transcription
01/06/2006 23:00	Play segment	Play file	Transcription
01/06/2006 20:00	Play segment	Play file	Transcription

Top 3 Topics

Topic 49 'California Politics' (probability 38.3%)

Topic Keywords:
california, southern, heat, temperatures, dollar, wave, weather, arnold, deaths, governor

Top 3 documents within topic:

25/07/2006 12:00	Play segment	Play file	Transcription
28/07/2006 05:00	Play segment	Play file	Transcription
25/06/2006 01:00	Play segment	Play file	Transcription

Topic 62 'Guard'

Topic Keywords:
guard, mexi, mexican, hu

Top 3 documents:

15/05/2006			
21/06/2006			
16/05/2006			

Topic 18 'State'

Topic Keywords:
state, gover, lawmakers,

Top 3 documents:

05/07/2006			
05/07/2006			
04/07/2006			

... california governor arnold's fortson agar inspected the california mexico border by helicopter wednesday to see ...

... the past days president bush asking california's governor for fifteen hundred more national guard troops to help patrol the mexican border but governor orville schwartz wicker denying the request saying...

Fig. 2. Two examples of the retrieved text for a query on 'schwarzenegger'.

castsearch.imm.dtu.dk

Modeling the generalizability of factorization

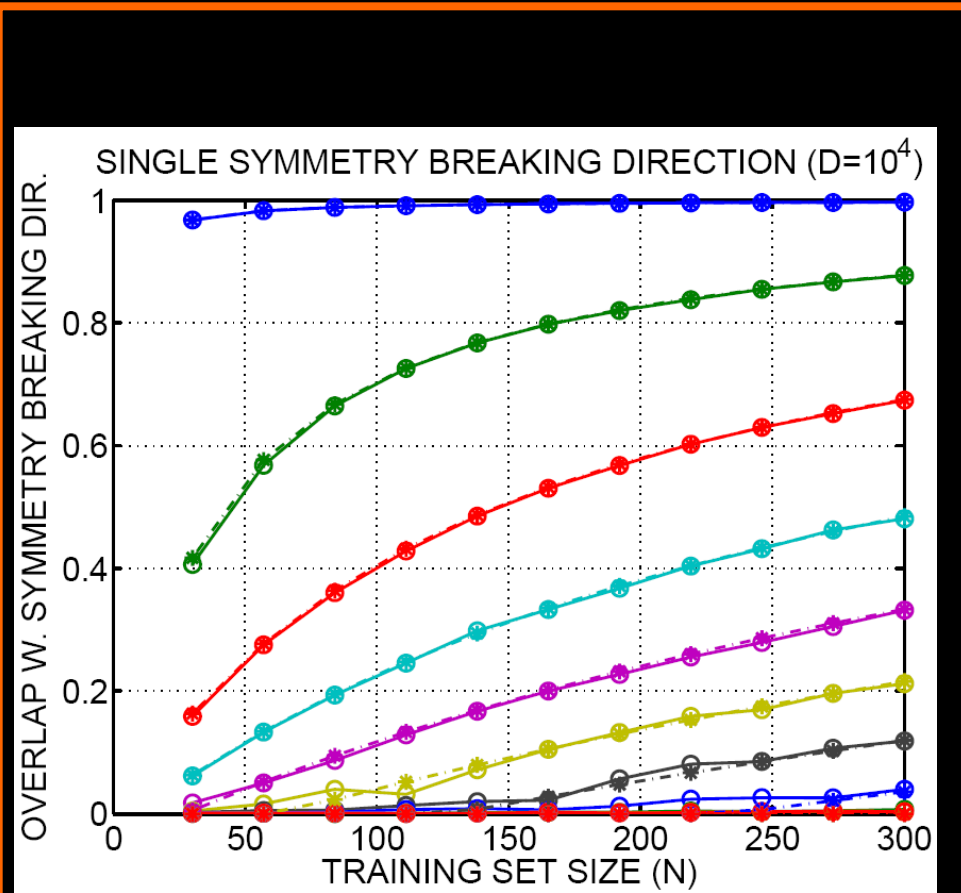


- Rich physics literature on "retarded" learning
- **Universality**
 - Generalization for a "single symmetry breaking direction" is a function of ratio of N/D and signal to noise S
 - For subspace models-- a bit more complicated -- depends on the component SNR's and eigenvalue separation
 - For a single direction, the mean squared overlap $R^2 = \langle (u_1^T u_0)^2 \rangle$ is computed for $N, D \rightarrow \infty$

$$R^2 = \begin{cases} (\alpha S^2 - 1) / S(1 + \alpha S) & \alpha > 1/S^2 \\ 0 & \alpha \leq 1/S^2 \end{cases}$$

$$\alpha = N/D \quad S = 1/\sigma^2 \quad N_c = D/S^2$$

Hoyle, Rattray: Phys Rev E 75 016101 (2007)



$N_c = (0.0001, 0.2, 2, 9, 27, 64, 128, 234, 400, 625)$

$\sigma = (0.01, 0.06, 0.12, 0.17, 0.23, 0.28, 0.34, 0.39, 0.45, 0.5)$

Quantifying subjectivity

"Affective computing"

- Affective computing is research in systems that can recognize, interpret, process, and simulate human emotion
- Emotions are omnipresent and extremely important to communication / opinion formation / intent reading etc
- Psychology of emotion is well developed but still far from complete...

7. The law of hedonic asymmetry: "Negative emotions last longer, positive emotions fade"



The Laws of Emotion

by Mico H. Frijda



Sentiment detection

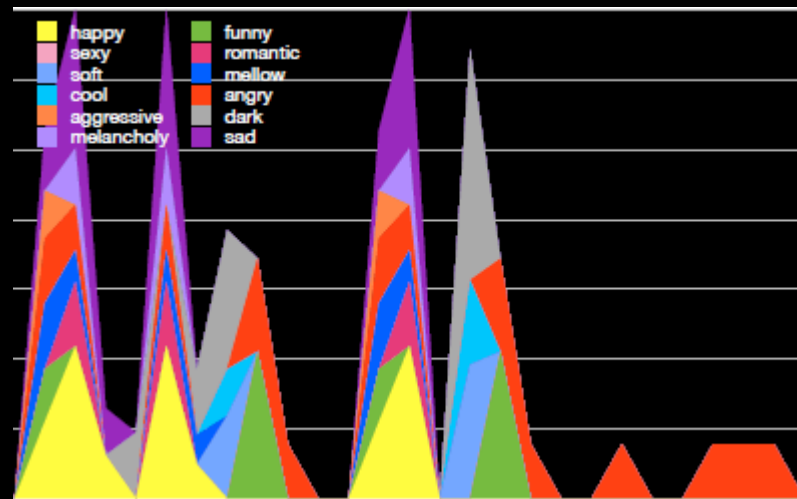
A step towards understanding subjectivity, opinion

Important to many services

Recommender (Amazon reviews)

Information navigation, e.g. navigating music

Oasis "Wonderwall"
Emotional content in the
song lyrics through time



M.K. Petersen: Modeling media as latent semantics
based on cognitive components (Ph.D. Thesis, DTU, 2010)

Text sentiment detection methods



WordNet

Use the linguistic structure (network of word relations) to compute the distance to "good" or "bad"

Supervised learning "seeing the stars"

Learn a predictive model based on labeled data, e.g. product reviews (text + # stars)

Curated word list

Design a list of domain relevant keywords with emotional annotation, e.g. the generic list ANEW

Defining virality

Virality lacks a formal definition,

- e.g., Wiktionary says: (*advertising, marketing*) "*The state or condition of being viral; tendency to spread by word of mouth.*"

We define virality statistically as the probability that a message is passed on in the network.

In Twitter this means retweeted, in other media different. Eg. in news papers one can *email* articles to friends (NYT)

Analysis of retweeting (Suh et al, 2010)



TABLE I. TWEET VARIABLES

<i>URL</i>	# of URLs in a tweet
<i>Hashtag</i>	# of hashtags in a tweet
<i>Mention</i>	# of usernames specified in a tweet excluding ones used for making a retweet (e.g. via @username, RT: @username)
<i>Follower</i>	# of users who follows the author of a tweet
<i>Followee (Friend)</i>	# of friends that the author is following
<i>Day</i>	# of days since the author of a tweet created the Twitter account
<i>Status</i>	# of tweets made by the user since the creation of the account
<i>Favorite</i>	# of favorited tweets by a user
<i>Retweet</i>	# of retweets recorded for a given tweet

TABLE IV. GENERALIZED LINEAR MODEL

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-4.42000	0.146400	-30.19	0.0000*
Days	0.00122	0.000296	4.12	0.0000*
HashtagOrNot	1.32800	0.160300	8.28	0.0000*
MentionOrNot	-0.29490	0.166800	-1.77	0.0771
URLOrNot	0.76360	0.150900	5.06	0.0000*
Followee	0.00006	0.000020	2.85	0.0043*
Follower	0.00002	0.000005	3.82	0.0001*
Status	-0.00002	0.000009	-1.71	0.0876
Favorite	-0.00004	0.000163	-0.26	0.7987

The Sentiment / Virality paradox

Seemingly conflicting observations on virality

i) Bad news is good news. Negative sentiment propagates in news media

(Galtung & Ruge, 1965)

ii) Most social networks including Twitter have a dominant positive sentiment

iii) emailing from NYT is dominated by positive sentiment articles (Berger & Milkman, 2010)

We will go to Twitter to find out!

DTU Informatics

Department of Informatics and Mathematical Modeling



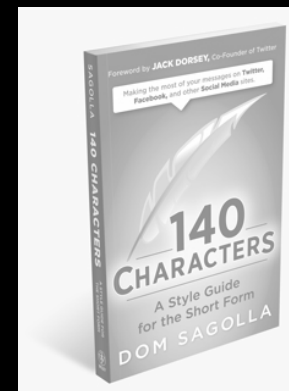
Twitter - a digital behavior lab

“Microblog”: Conceived as SMS/Texting Short Text Messaging (140 characters) for the Internet

Jack Dorsey, Biz Stone and a group of engineers defined Twitter in March 2006, and launched it in July 2006.

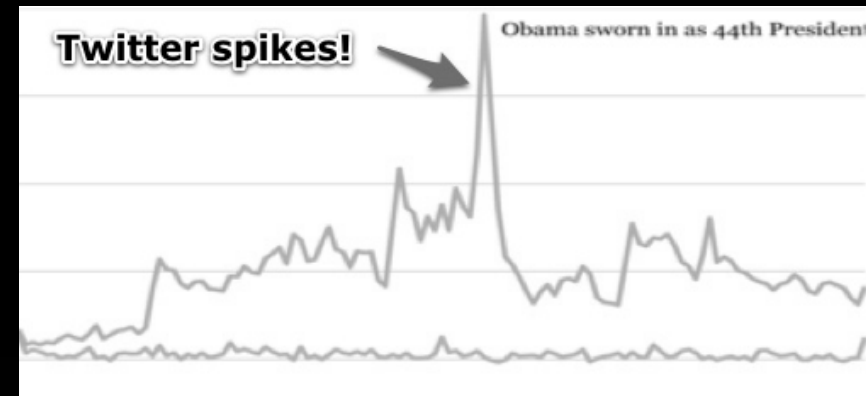
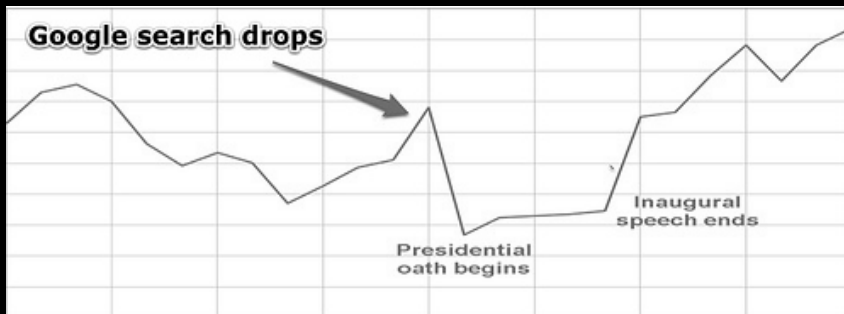
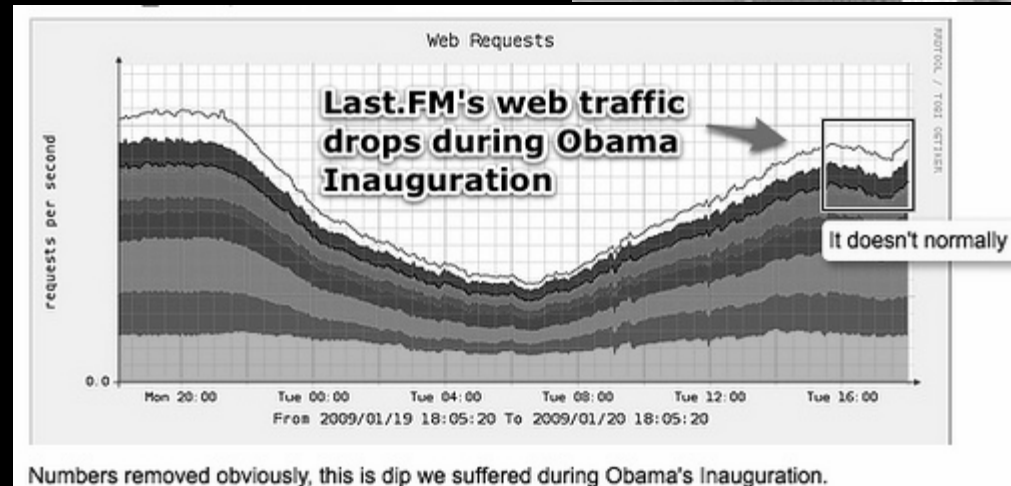
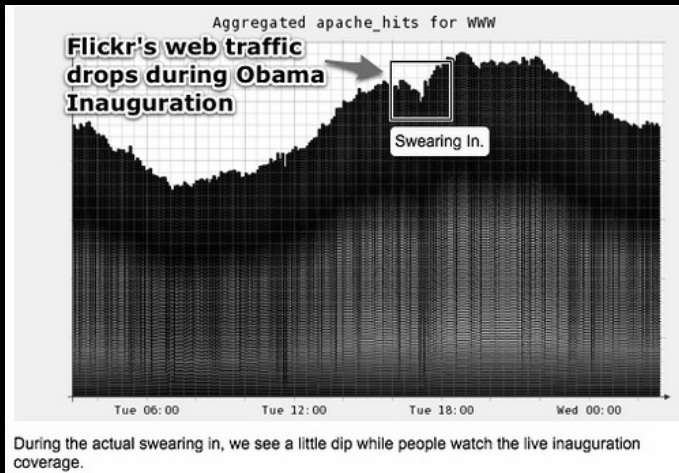
200 mill. users, 65 mill. messages pr day

Dom Sagolla: 140 Characters (Wiley, 2009)



Web of time...the blog pulse

Jesse Robbins O'reilly's radar Feb 08, 2009



DTU Informatics

Department of Informatics and Mathematical Modeling

Twitter is a news medium

Obama inauguration speech,

Egypt upraise tag #jan25

Tsunami, Japan phone system collapsed, but net still worked and Tweeting soared

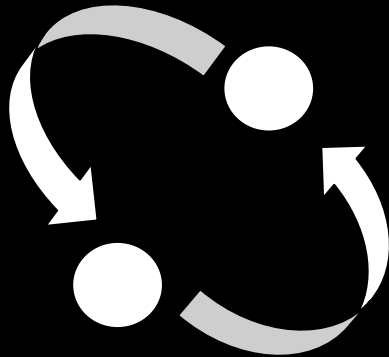


<http://mashable.com/2011/03/11/japan-tsunami/>

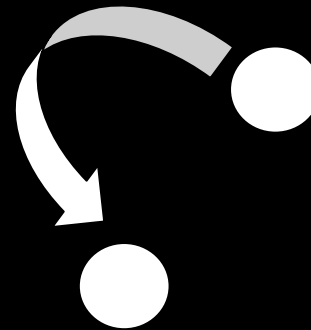
Twitter is a social medium

Twitter conceived as Internet SMS/texting – keeping friends updated

Symmetry of connections



Friendship is symmetric



News / interest graph may be asymmetric

So is Twitter a news medium or a social network?

Kwak et al. (2010)

"...Twitter shows a low level of reciprocity; 77.9% of user pairs with any link between them are connected one-way, and only 22.1% have reciprocal relationship between them.

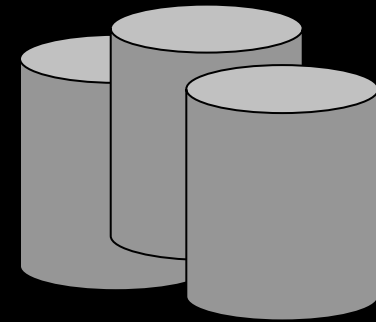
"...Previous studies have reported much higher reciprocity on other social networking services: 68% on Flickr [4] and 84% on Yahoo! 360.

"...Moreover, 67.6% of users are not followed by any of their followings in Twitter. We conjecture that for these users Twitter is rather a source of information than a social networking"

Understanding virality and sentiment in Twitter

Our research questions

- Q1: How accurately can text be characterized as 'news'?
- Q2: How big a fraction of Twitter is news?
- Q3: If Twitter is a news medium, does negative sentiment influence virality?
- Q4: Does sentiment influence retweet probability differentially in news and social messages?



Databases for this study

Brown corpus

a general text corpus with a known mixture of news/nonnews documents. The corpus consist of 47134 sentences.

RANDOM Twitter sample

348862 tweets collected September 9-14, 2010. The Tweets were randomly sampled following the 'Spritzer' protocol.

COP15 Twitter sample

complete set of tweets for a specific news event COP15 2009 UN Climate Change Conference, Denmark Dec 7-18. 207782 tweets downloaded December 1-31 with the term/tag "cop15".

News detection (Q1: Browne corpus)

$$\begin{aligned} p(\text{news}|\mathbf{w}) &= \frac{p(\text{news})p(\mathbf{w}|\text{news})}{p(\text{news})p(\mathbf{w}|\text{news}) + p(\neg\text{news})p(\mathbf{w}|\neg\text{news})} \\ &= \left(1 + \frac{p(\neg\text{news}) \prod_{d=1}^D p(w_d|\neg\text{news})}{p(\text{news}) \prod_{d=1}^D p(w_d|\text{news})} \right)^{-1} \end{aligned}$$

We posit a simple model (Naive Bayes) to estimate the probability of carrying the news label, given the bag of words features (D=10000). The accuracy on test data is 84%

News detection (Q2: Twitter news?)

We apply the trained NB classifier to the two
Twitter samples

A tweet is declared news if $p(\text{news} \mid w) > 0.5$

RANDOM data

Rate of news: 22.3%

COP15 data

Rate of news: 30.3%

Sentiment detection in Twitter (Q3+Q4 How does sentiment affect retweeting?)

Now we are interested in estimating the probability of retweet and to infer the importance of the features.

We use the generalized linear model as in Suh et al., it provides standard scores for deletion of individual features (saliency)

$$p(RT|\mathbf{f}) = \left(1 + e^{-\sum_{i=0}^F \beta_i f_i}\right)^{-1}$$

Sentiment detection in Twitter (Q3+Q4 How does sentiment affect retweeting?)

GLM Features: Negative sentiment (*news), mention, url, hash-tag

RANDOM data set

There is no significant dependency on presence of negative sentiment.

Negative AND news, is retweeted more.

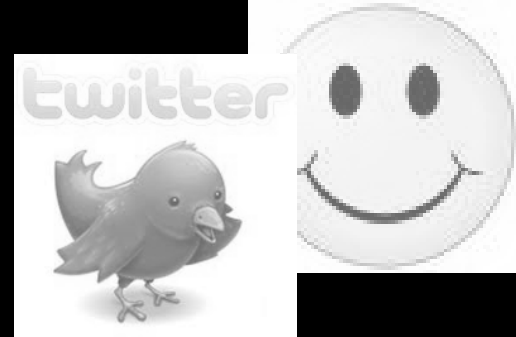
Confining to tweets that have a non-zero sentiment retweeting strongly suppressed by negative sentiment

COP15 data set (a news event)

For all tweets in this sample there significant effect that negative sentiment promotes retweeting. Also seen in the sentimental tweets.

Sentiment detection in Twitter (Q3+Q4 How does sentiment affect retweeting?)

GLM		COP15		RANDOM	
		All	English	All	English
RAN	N	147,041	136,262	335,236	106,719
E	Rate News	0.303	0.305	0.226	0.233
C	t(Negative)	4.889	4.649	2.775	-0.024
t	t(Negative_newsness)	2.275	1.471	-6.019	3.904
S	Tweets with Arousal > 0				
S	N	44,611	42,087	53,473	51,929
COP	t(Negative)	3.276	2.372	-9.725	-9.374
E	t(Negative_newsness)	1.125	0.180	1.179	1.239



Conclusions

- We may train computational models to partially understand social media behaviors
- Twitter is both a social and a news medium
- Virality can be defined through the retweet probability (probability of message being passed on).
- Likelihood of retweet is increased by negative content in newsy posts, in line with general news media theory
- In general propagation in social media is enhanced by positive content.



DTU Informatics

Department of Informatics and Mathematical Modeling

Contact

Prof. Lars Kai Hansen

Department of Signal Theory and Communications,
Universidad Carlos III de Madrid, 28911 Leganes, Spain

Permanent address:

DTU Informatics

Richard Petersens Plads B321

DK-2800 Kgs. Lyngby

lkh@imm.dtu.dk

www.imm.dtu.dk/~lkh

Tel.: 4525 3889

DTU Informatics

Department of Informatics and Mathematical Modeling
