

# Surprise!

*Information Surprise  
or how to discover data*

Oleksandr Pryymak /Sasha/  
@opryymak  
opryymak@gmail.com



*What is a ‘surprise’?*

# surprise

[countable] an event, a piece of news, etc. that is unexpected or that happens suddenly

SYNONYMS: *shock, ... , eye-opener*

[uncountable, countable] a feeling caused by something happening suddenly or unexpectedly

SYNONYMS: *astonishment, ...*

(Oxford Advanced Learner's Dictionary)

*Define Surprise!*

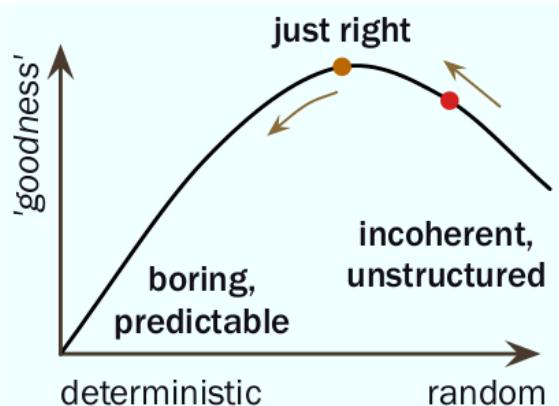


measured in  
**wows**

*Quantify Surprise!*

# Quantify

Need for measuring any type of content: **complexity measure**  
Complex is not randomness!

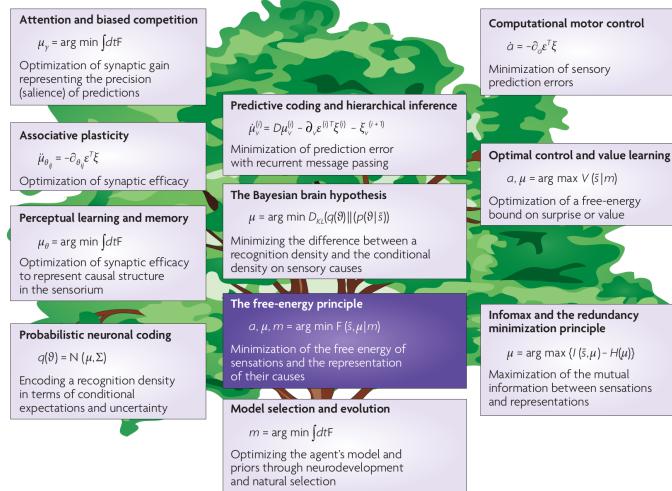


## Measures of complexity

1. Subjective rating
2. #Distinct elements
3. #Dimension
4. #Control parameters
5. Minimal description
6. Information content
7. Minimal generator
8. Minimum energy

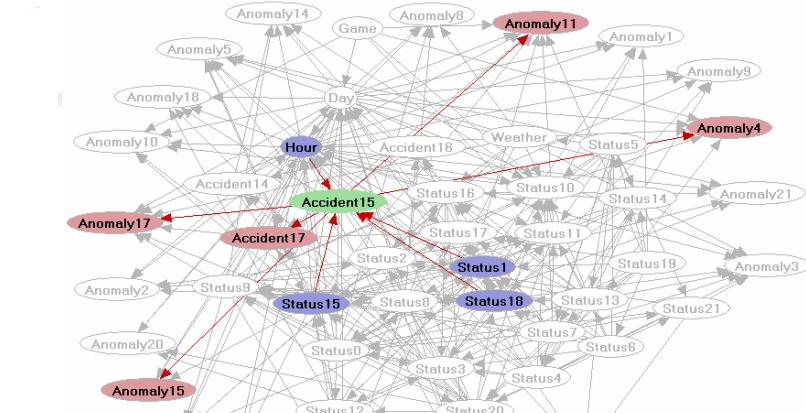
# Surprise Quants in academia

## Neuro/Cognitive Science *How do we perceive information?*



Friston, K. (2010). *The free-energy principle: a unified brain theory?*. *Nature Reviews Neuroscience*, 11(2), 127-138.

## Machine Learning *How to measure differences?*

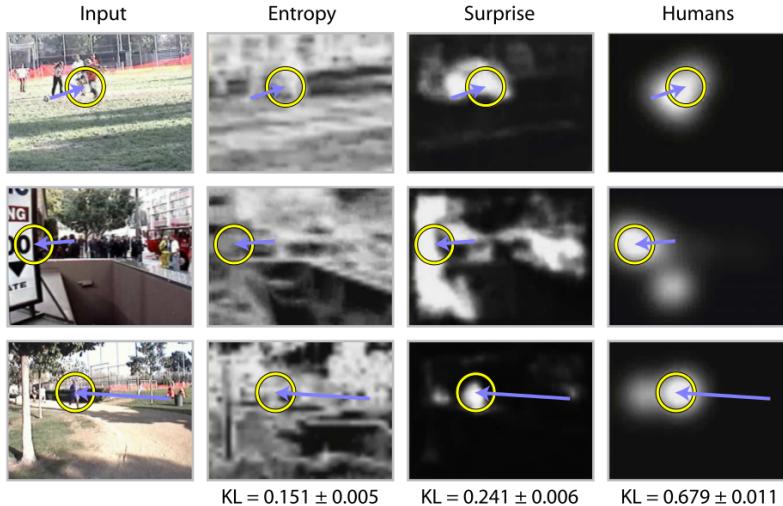
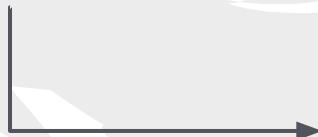


... machine that constantly tells you what you already know is just irritating. So software alerts users only to surprises...

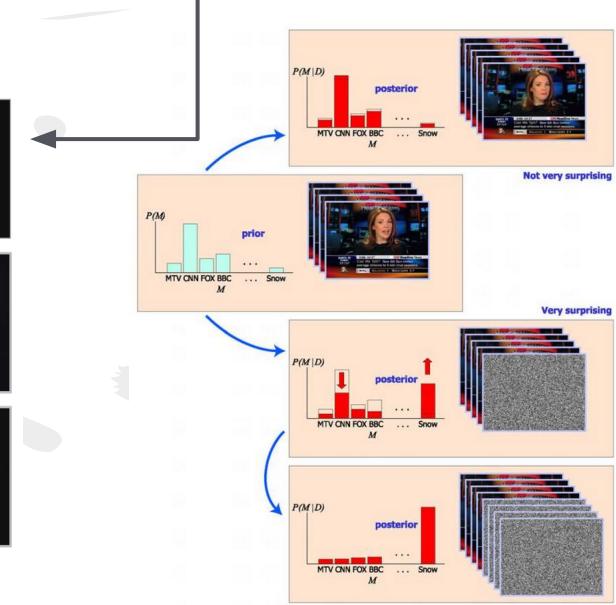
Horvitz, E., Apacible, J., Sarin, R., & Liao, L. *Prediction, Expectation, and Surprise: Methods, Designs, and Study of a Deployed Traffic Forecasting Service*.

# Surprise Quants in academia

Neuro/Cognitive Science



Machine Learning



Itti, L., & Baldi, P. F. (2005). *Bayesian surprise attracts human attention*. In *Advances in neural information processing systems* (pp. 547-554).

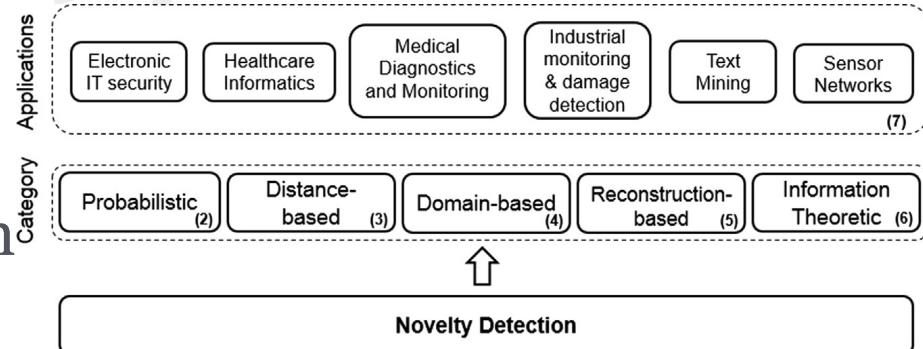
# Typical ML applications

## Unsupervised Learning

1. Decision trees (inf. gain)
2. MaxEnt principle
3. ...

Specifically after ‘surprise’:

4. One-class classification
5. Anomaly detection
6. Novelty measure



Pimentel, M. A., Clifton, D. A., Clifton, L., & Tarassenko, L. (2014).  
A review of novelty detection. *Signal Processing*, 99, 215–249.

# Quantify surprisal /self-information/

The surprise /information/ we get from observing the occurrence of an event having probability p.

Axioms:

$$0 < p \leq 1$$

$$I(p) \geq 0$$

$$I(1) = 0$$

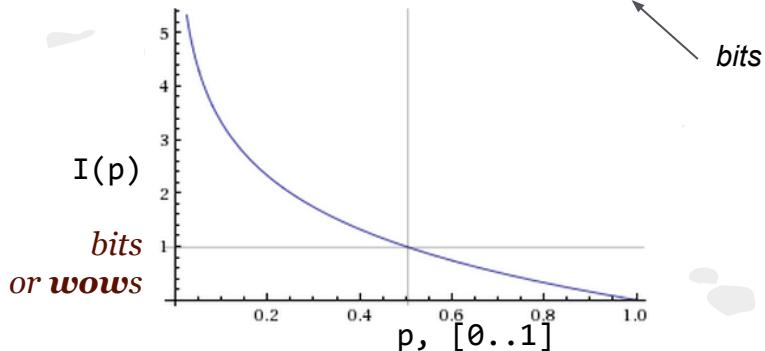
$$I(p_1 * p_2) = I(p_1) + I(p_2) \text{ //independent}$$

Derive:

$$I(p^2) = I(p * p) = I(p) + I(p) = 2 * I(p)$$

$$I(p^a) = a * I(p)$$

Surprisal /self-information/:

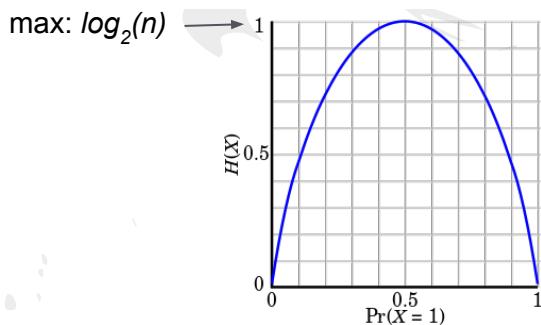
$$I(p) = -\log_2(p) = \log_2(1/p)$$


Flipping a fair coin provides 1bit of new information.

# Quantify ‘knowledge’ /entropy/

The Shannon *entropy* is the expected value of the *self-information*.

$$H(P) = \langle I(p) \rangle = \sum_{i=1}^n p_i * \log_2(1/p_i)$$



Entropy of a Bernoulli trial  
 $X \in \{0,1\}$

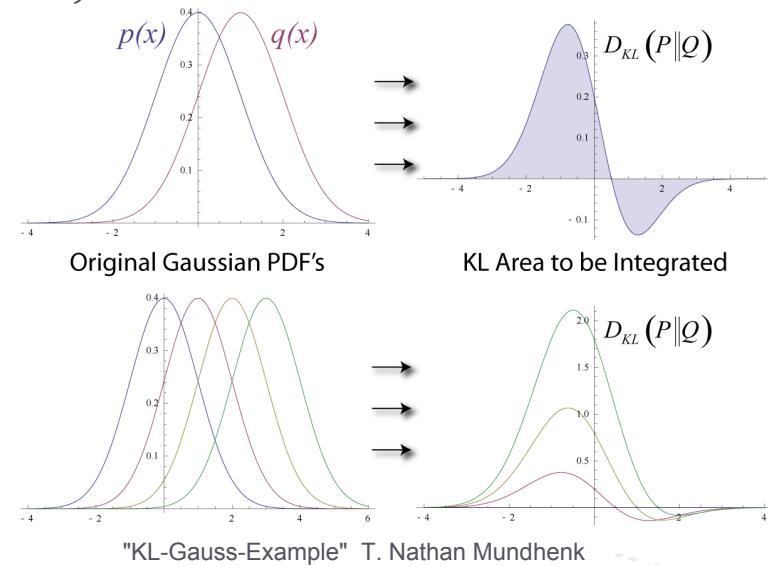
Notes:

1. The maximum entropy distribution is the least informative.
2. The *statistical mechanics* and the *information* entropy are principally the same.

# Quantify ‘discovery’ /information gain/

The **Kullback–Leibler divergence** /relative entropy, information gain/:  
is a measure of the information lost when Q is used to approximate P  
(measures the expected number of extra bits required to recode)

$$D_{KL}(P||Q) = \sum_{i=1}^n p_i \log_2 \frac{p_i}{q_i}$$



Assymetric: not a true measure

# Quantify ‘discovery surprise’

Symmetric KL Distances:

$$D_{KL2}(P, Q) = \sum_{i=1}^n (p_i - q_i) \log_2 \frac{p_i}{q_i}$$

$$D_{KL1}(P, Q) = D_{KL}(P||Q) + D_{KL}(Q||P)$$

$$D_{KL3}(P, Q) = \frac{1}{2} \left[ D_{KL}\left(P\middle|\frac{P+Q}{2}\right) + D_{KL}\left(Q\middle|\frac{P+Q}{2}\right) \right]$$

$$D_{KLA}(P, Q) = \max(D_{KL}(P||Q), D_{KL}(Q||P))$$

All result in the same performance:

Pinto, D., Benedí, J. M., & Rosso, P. (2007). Clustering narrow-domain short texts by using the Kullback-Leibler distance. In *Computational Linguistics and Intelligent Text Processing*

# Calculating KLD

Data sparseness problem: often  $p_i=0 \rightarrow \text{KLD}(P||Q)=\infty$

Solutions:

- drop components from calculations
- smoothing:

$$p_i = \begin{cases} \beta * p_i, & \text{if } p_i > 0 \\ \epsilon = 0.0001 & \text{otherwise} \end{cases}$$

$$\beta = 1 - \sum_{\forall p_i=0} \epsilon$$



BBC News (World)

Intense clashes. Riot police falling back.  
Petrol bombs, tear gas, a lot of smoke.  
#Kiev #Ukraine - via @\_DuncanC

RETWEETS 362 FAVORITES 74

7:35 AM - 20 Feb 2014

Reply to @BBCWorld @\_DuncanC

listyo budi santoso @listyobudi · 20 Feb 2014  
@BBCWorld @\_DuncanC Oh meenn...

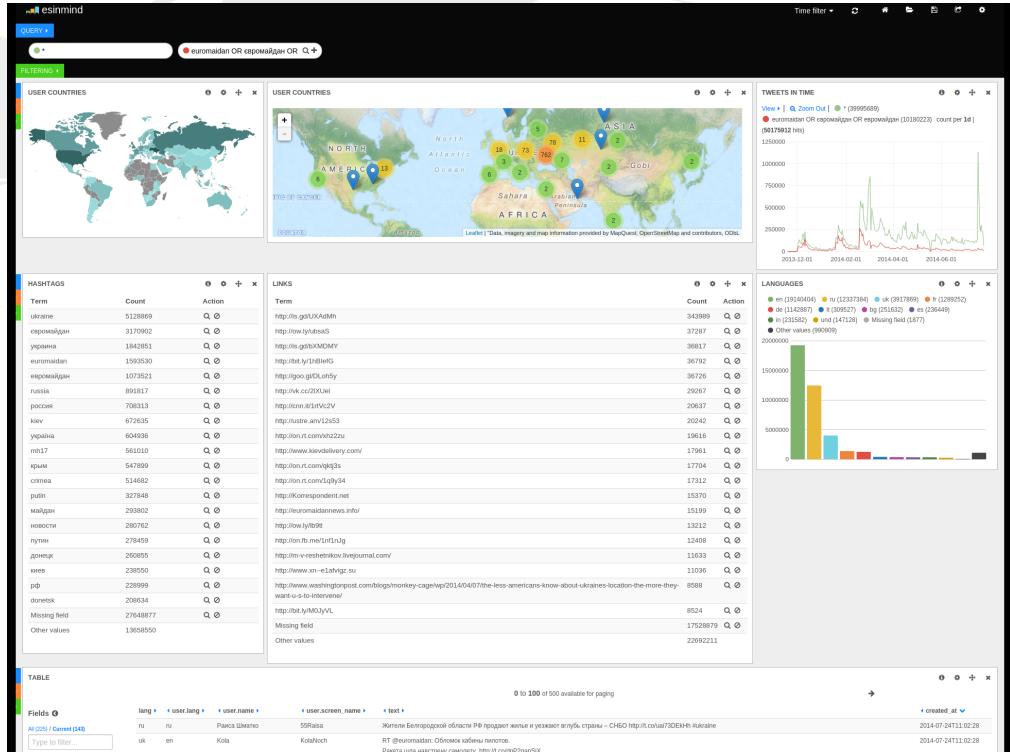


*Surprise in Twitter feed*

# Search engines

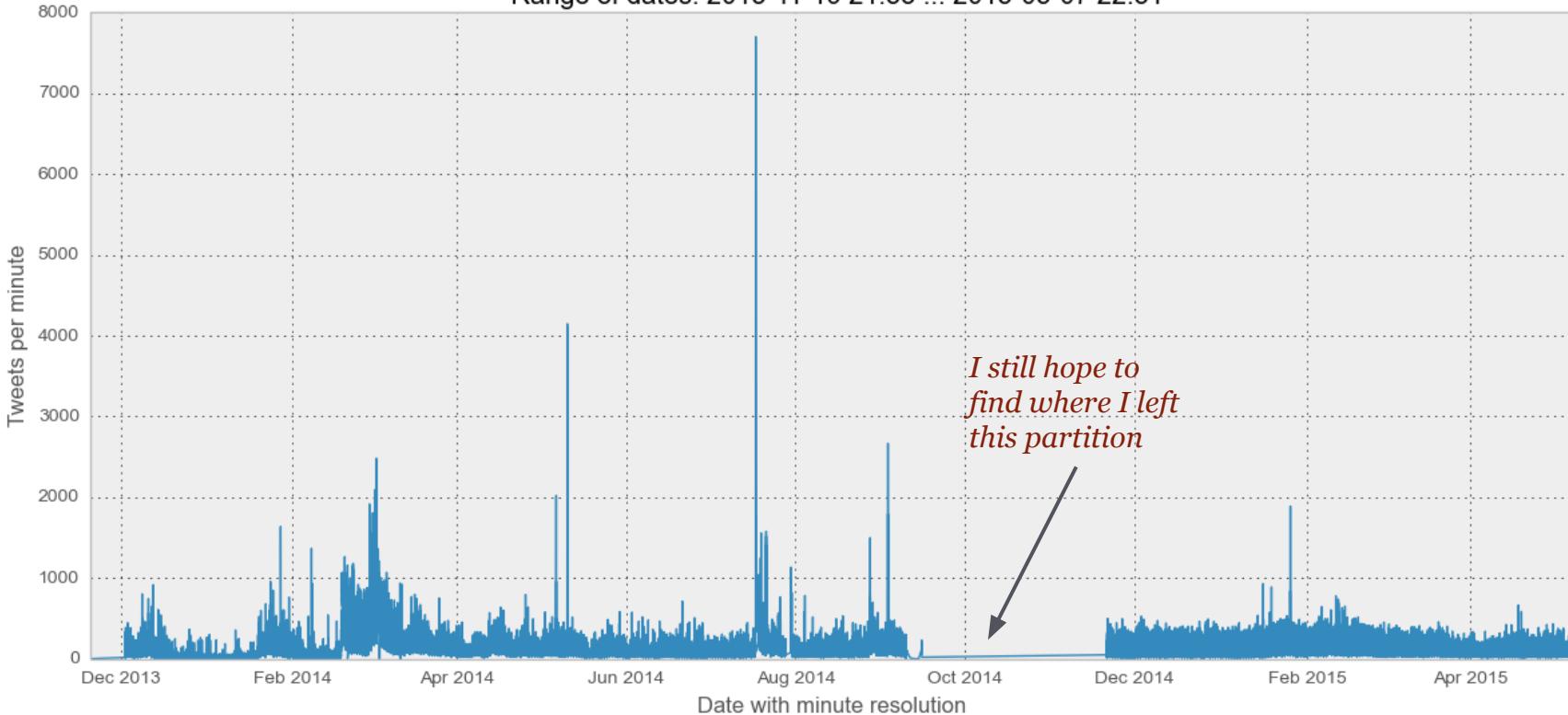
Elasticsearch +  
Kibana =  
faceted data exploration

(tweets data)



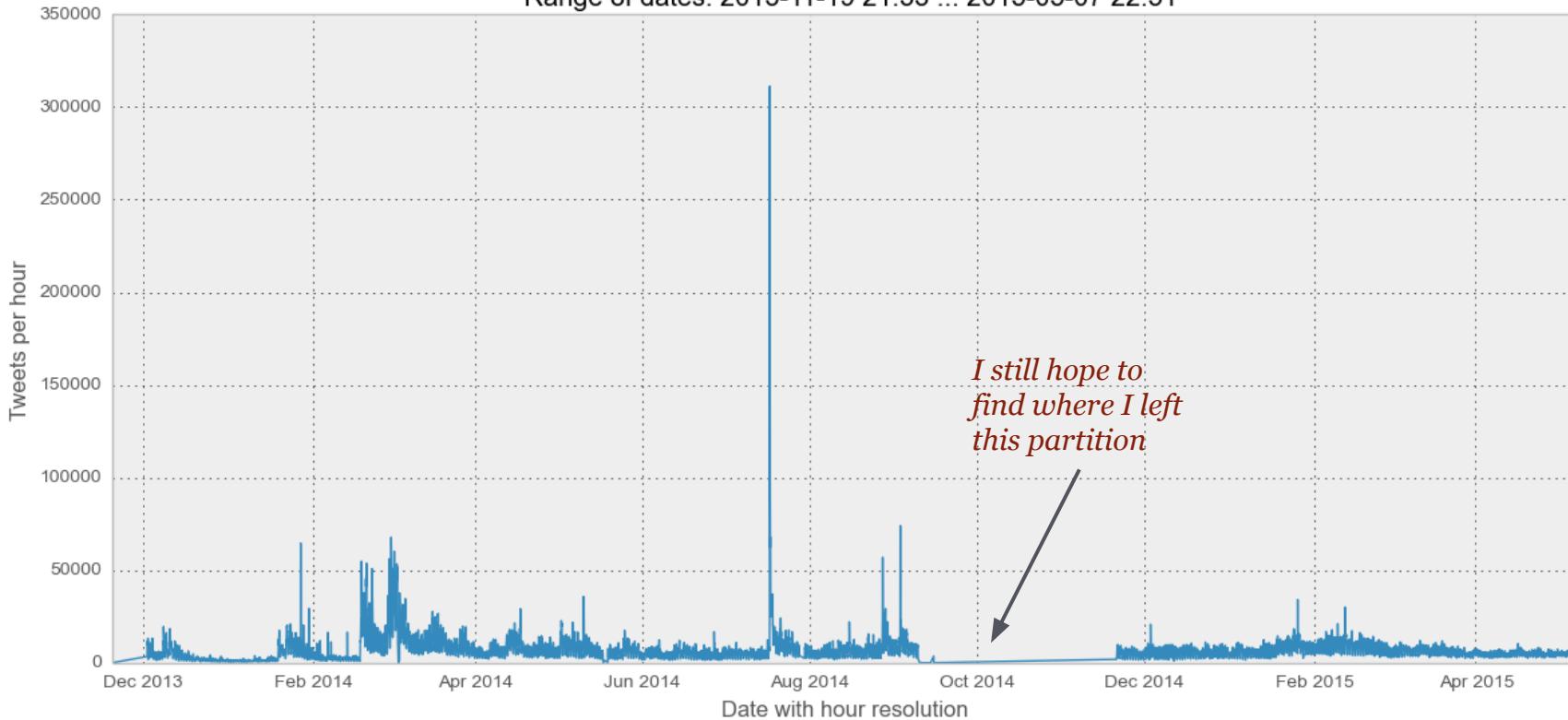
# Whole dataset

Tweets per minute (total 75013738 tweets)  
Range of dates: 2013-11-19 21:53 ... 2015-05-07 22:31



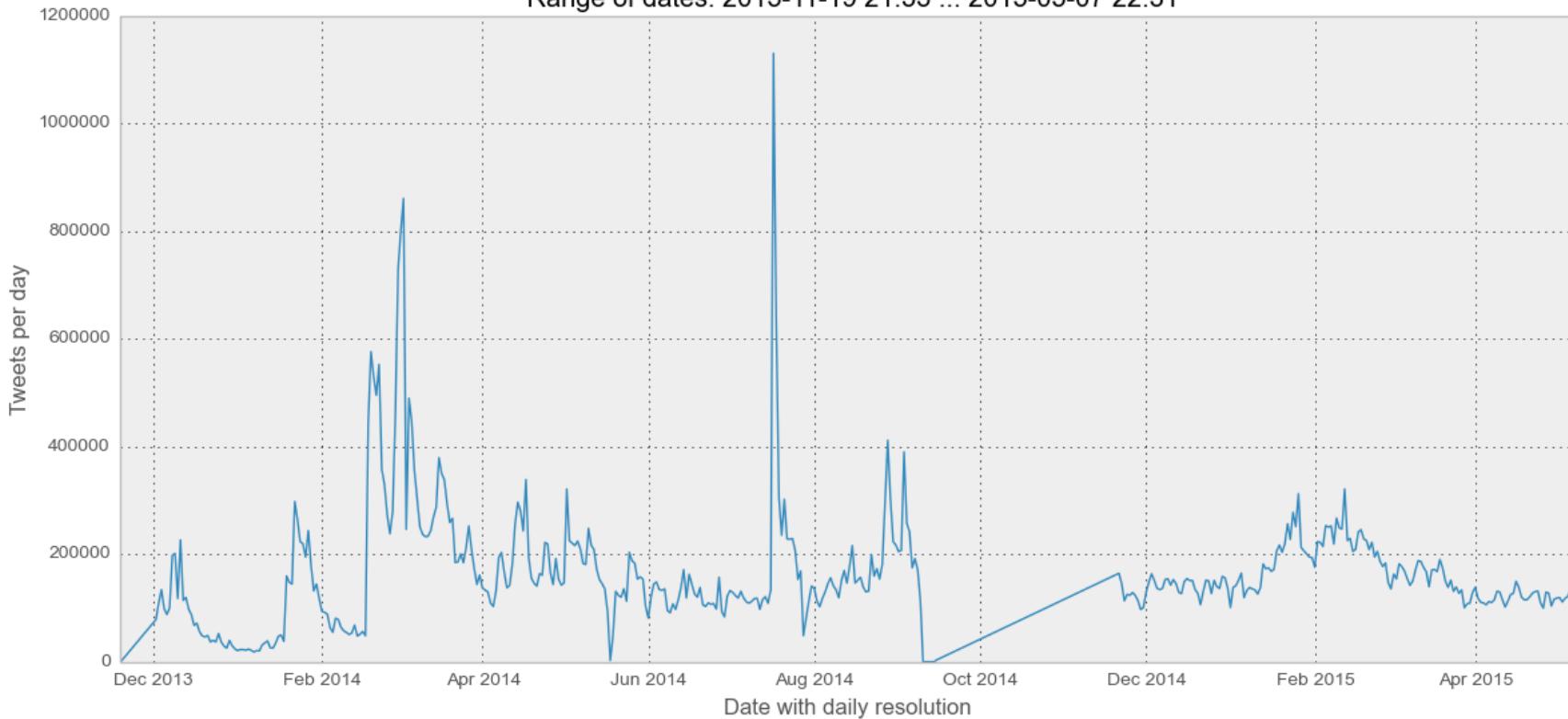
# Whole dataset

Tweets per hour (total 75013738 tweets)  
Range of dates: 2013-11-19 21:53 ... 2015-05-07 22:31

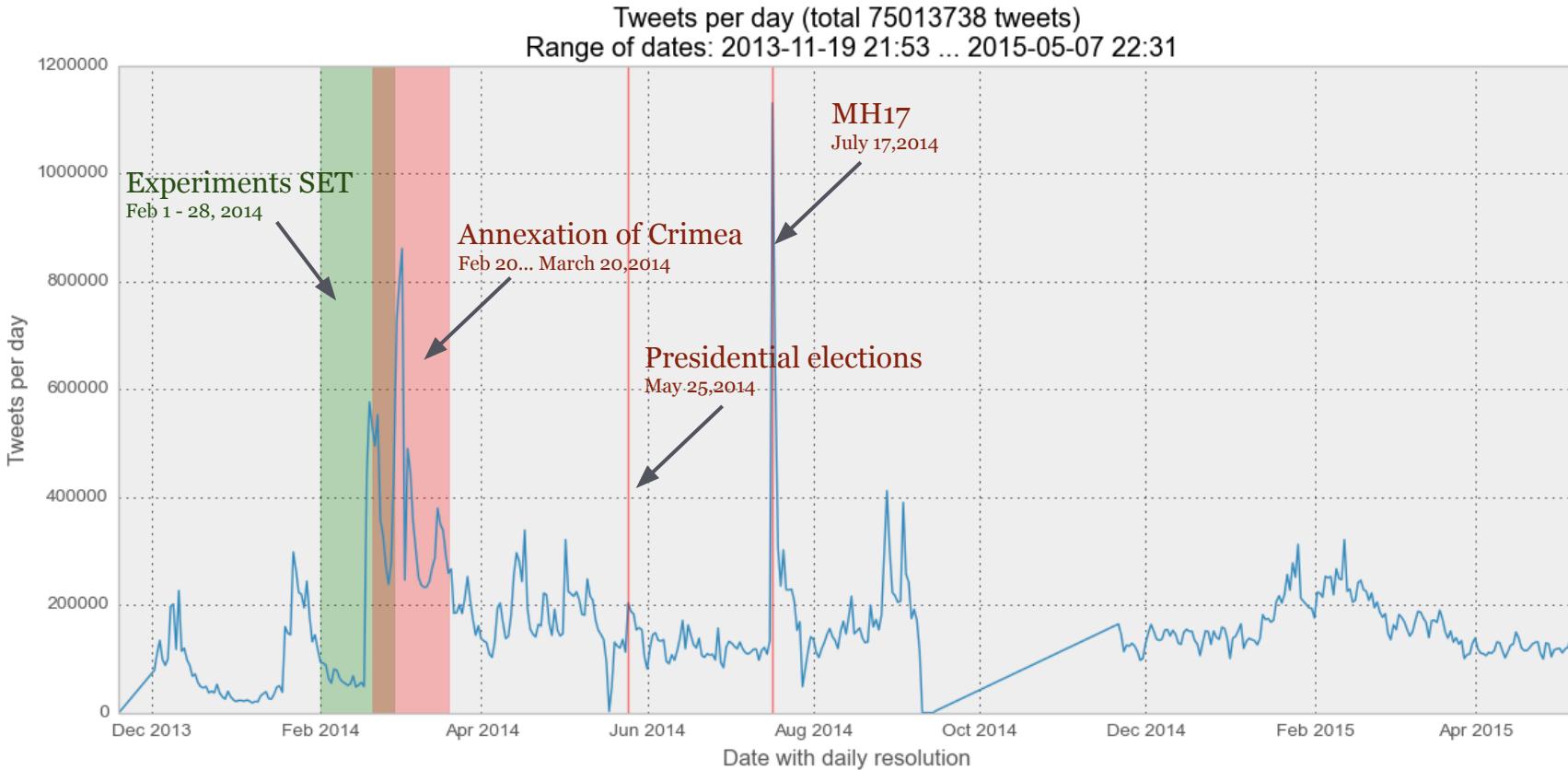


# Whole dataset

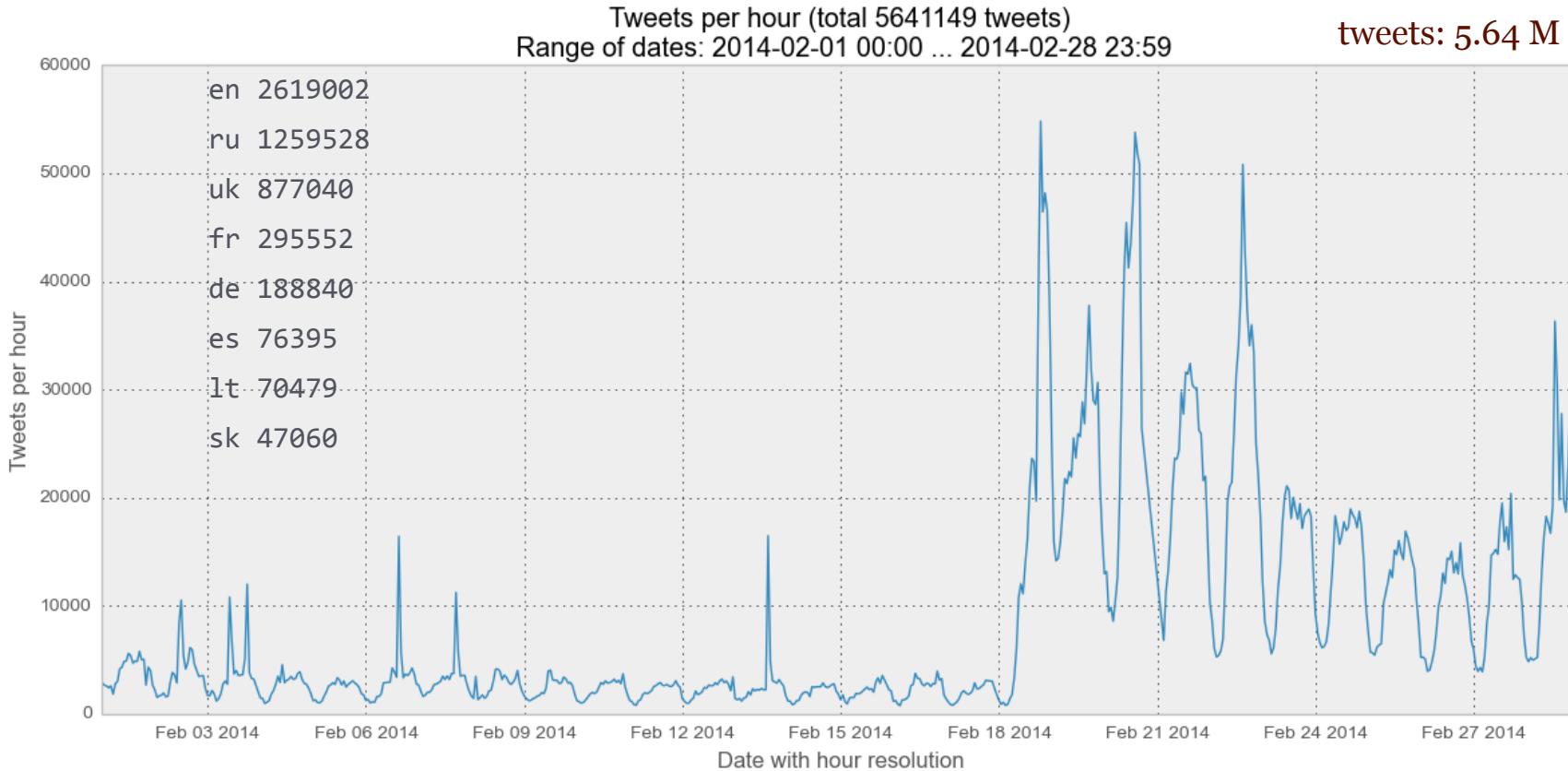
Tweets per day (total 75013738 tweets)  
Range of dates: 2013-11-19 21:53 ... 2015-05-07 22:31



# Whole dataset

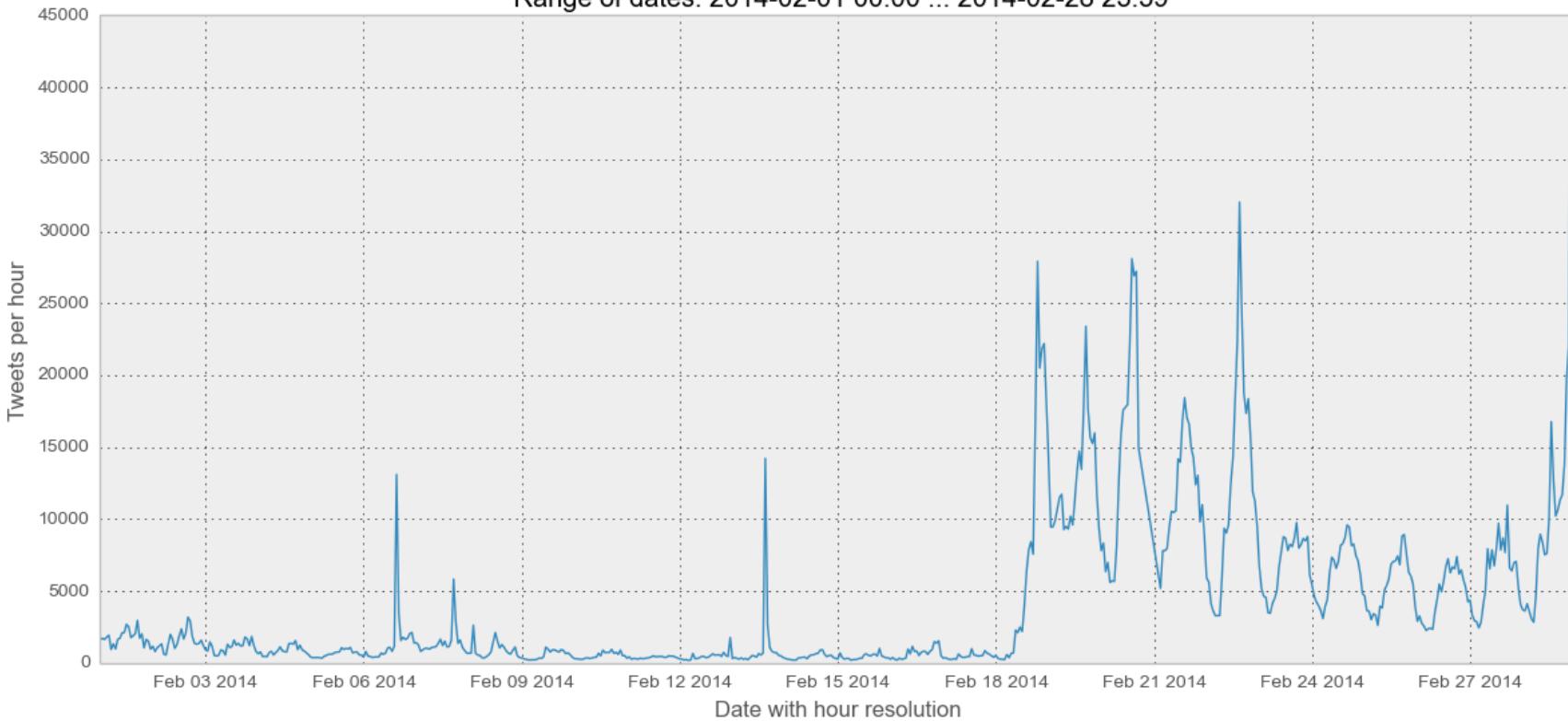


# Experiment dataset: Feb2014



# Experiment dataset: English

ENG Tweets per hour (total 2619002 tweets)  
Range of dates: 2014-02-01 00:00 ... 2014-02-28 23:59



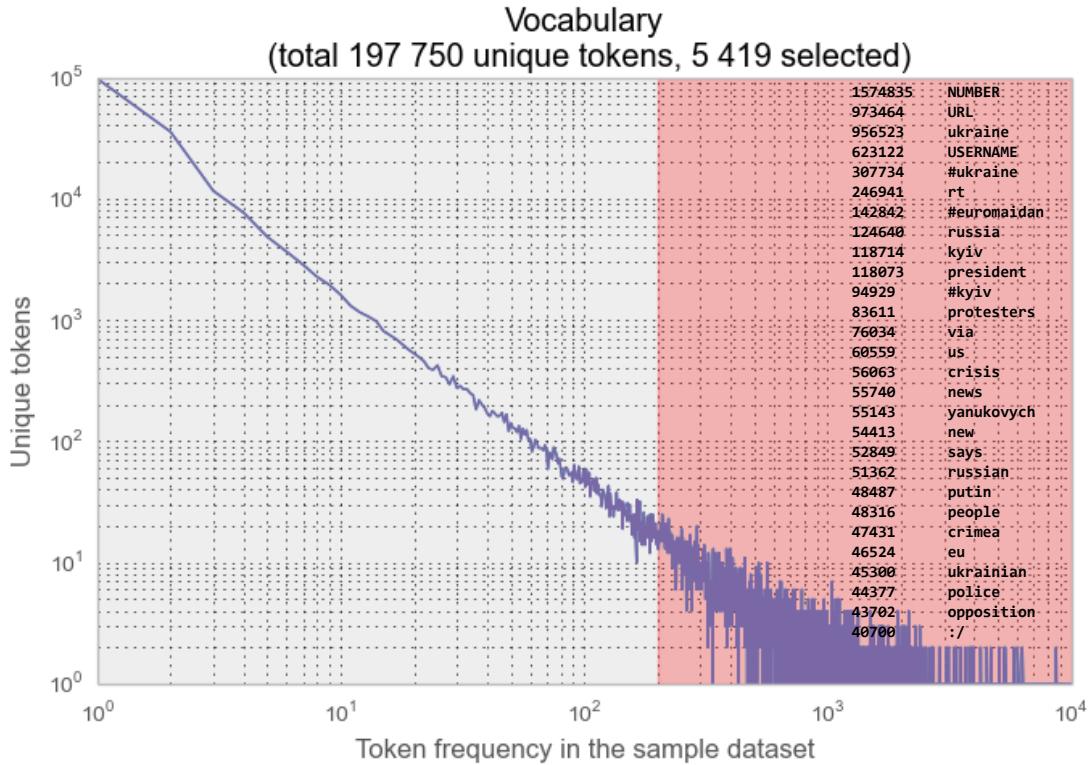
# Topic modeling

Tweets are short, but prominent events are widely discussed.

Document is a time slot.

Model:

- bag of words
- freq. threshold > 200 tweets
- term frequency (naive)
- tokenizer: [https://github.com/jaredks/tweet\\_tokenize](https://github.com/jaredks/tweet_tokenize)  
+ a few touches



# Event Detection Problem

## Outliers detection:

- surprising
  - tweet rate

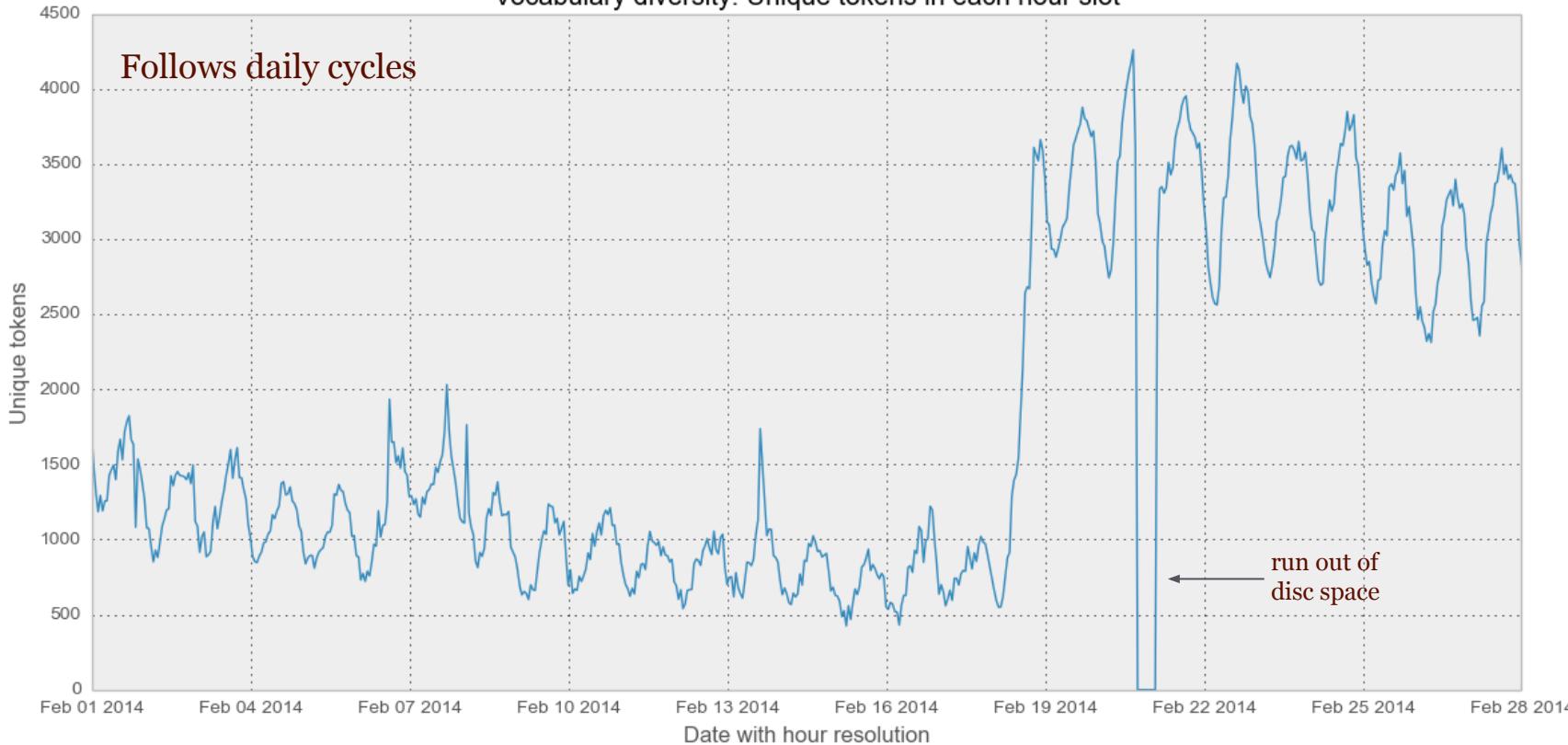
# Compare against:

historic ‘ground truth’ ->

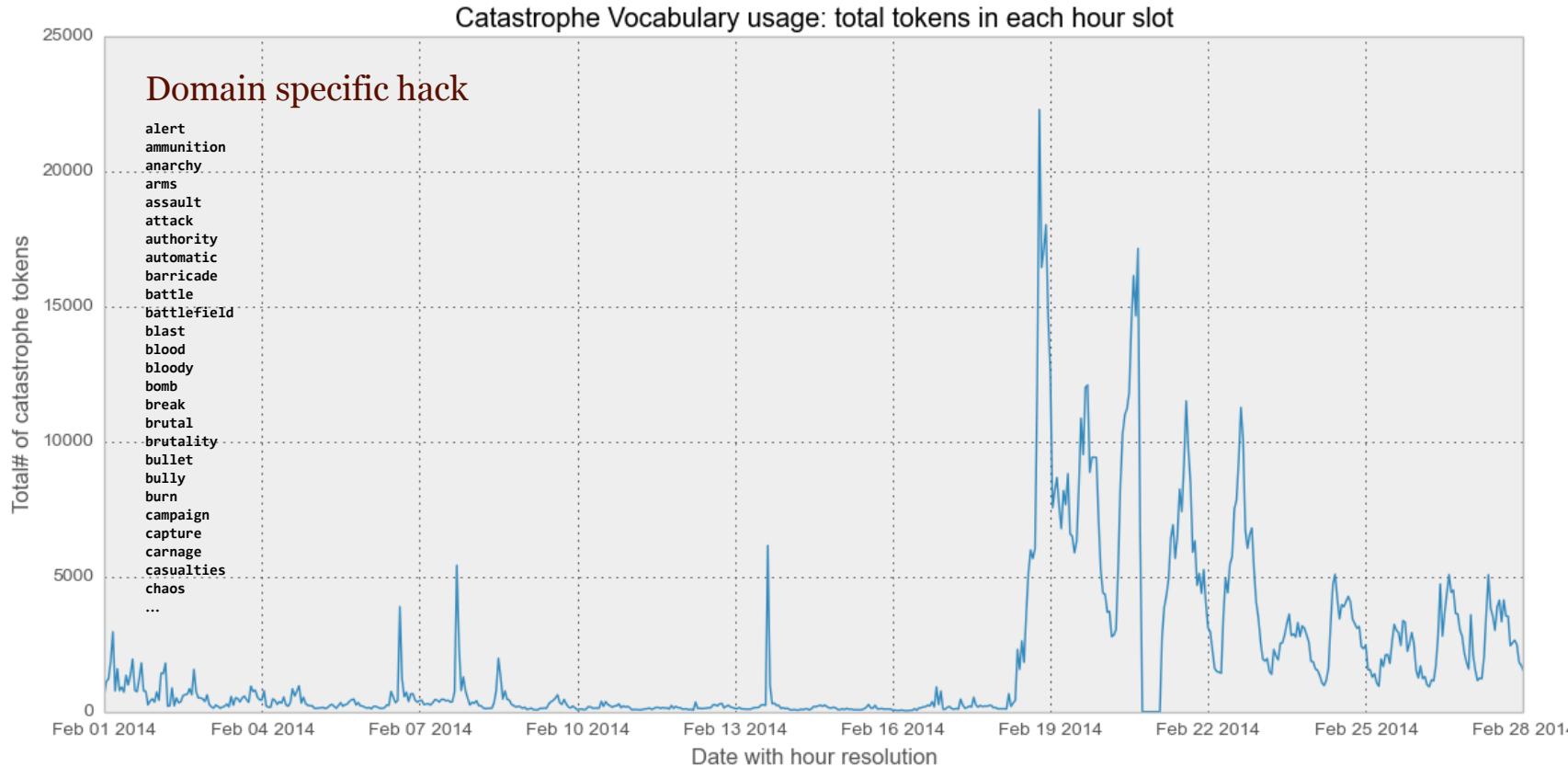


# Vocabulary diversity

Vocabulary diversity: Unique tokens in each hour slot



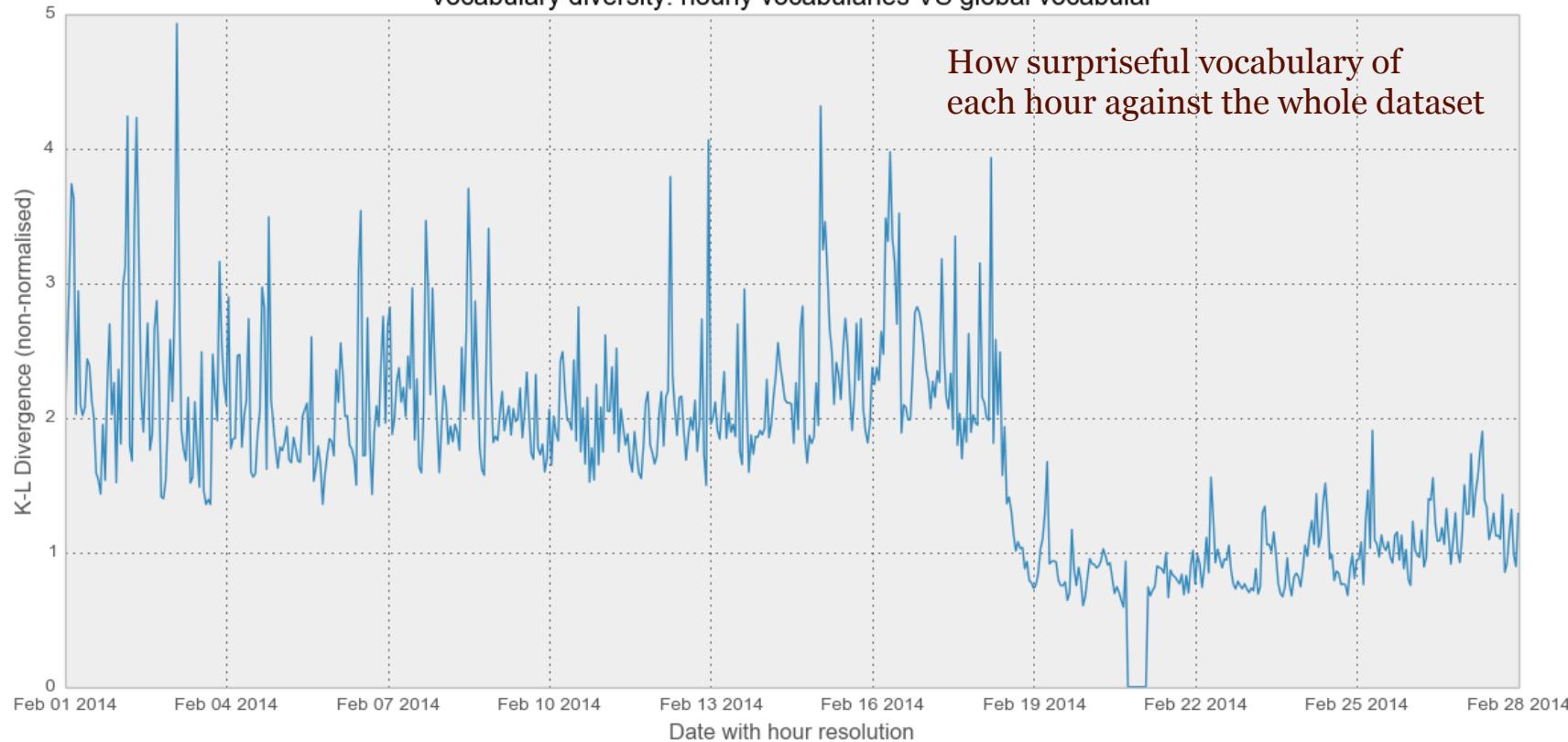
# Vocabulary: catastrophe



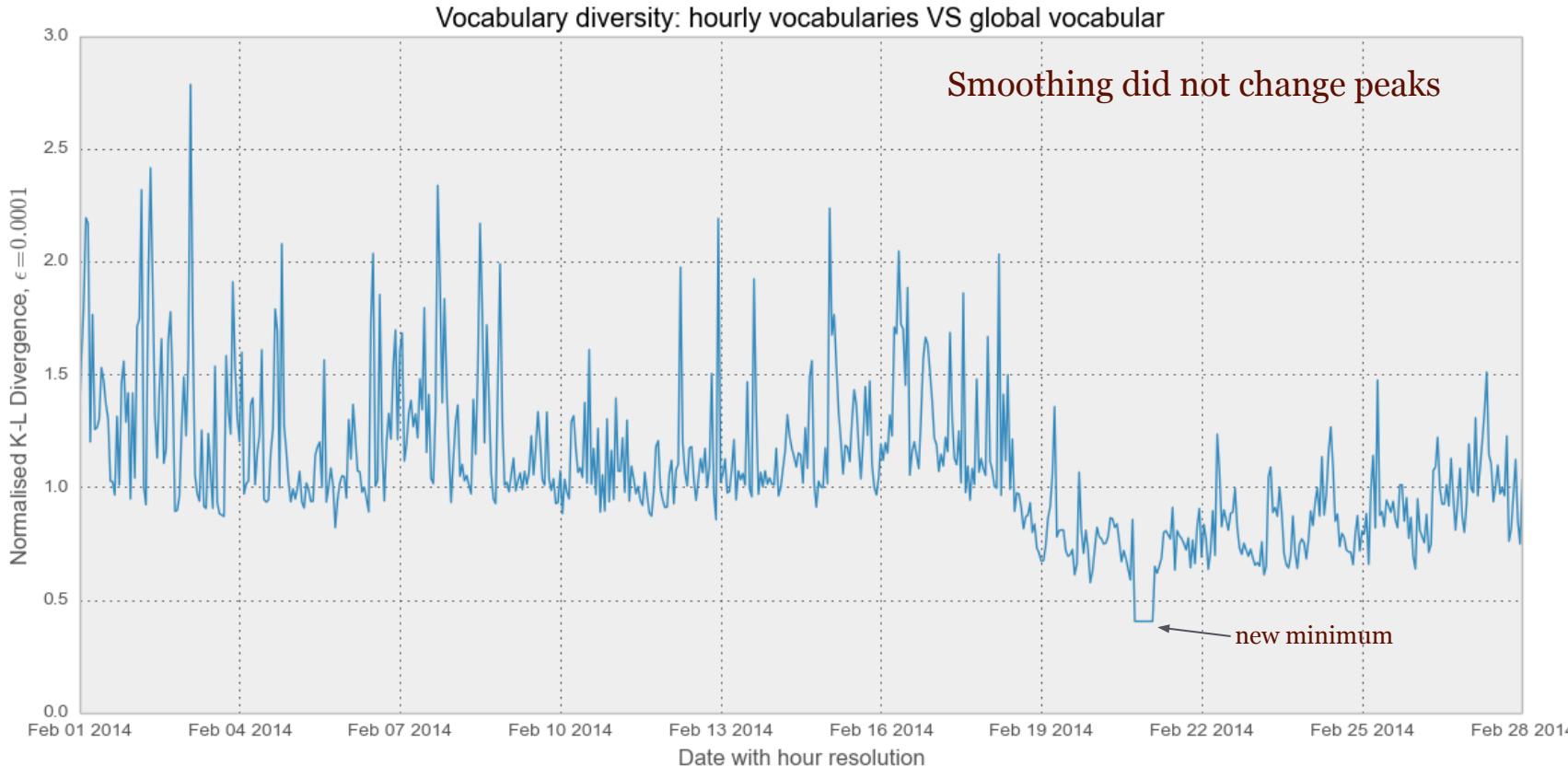
# Vocabulary slots: KLD

Vocabulary diversity: hourly vocabularies VS global vocabulary

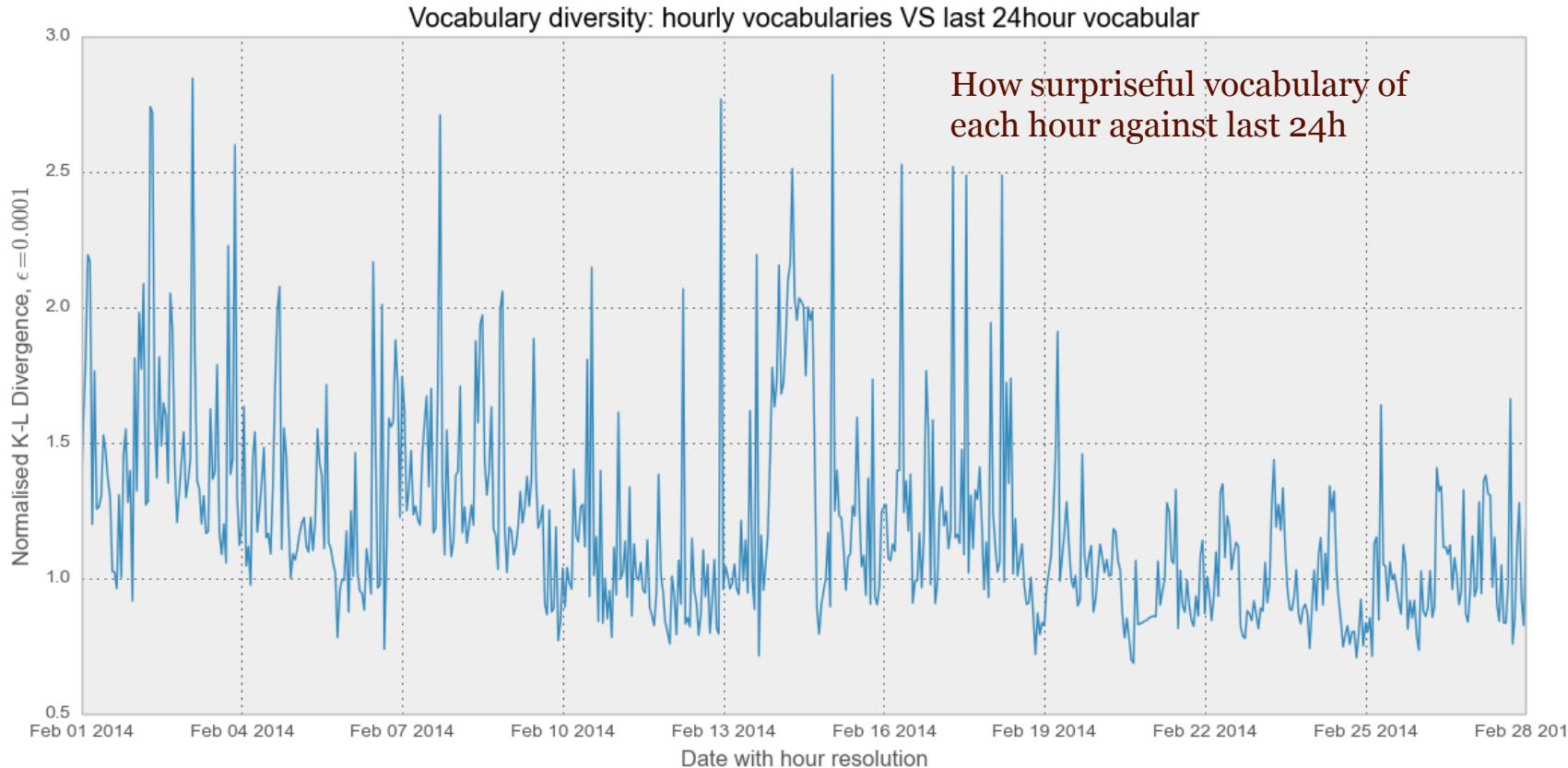
How surprised vocabulary of each hour against the whole dataset



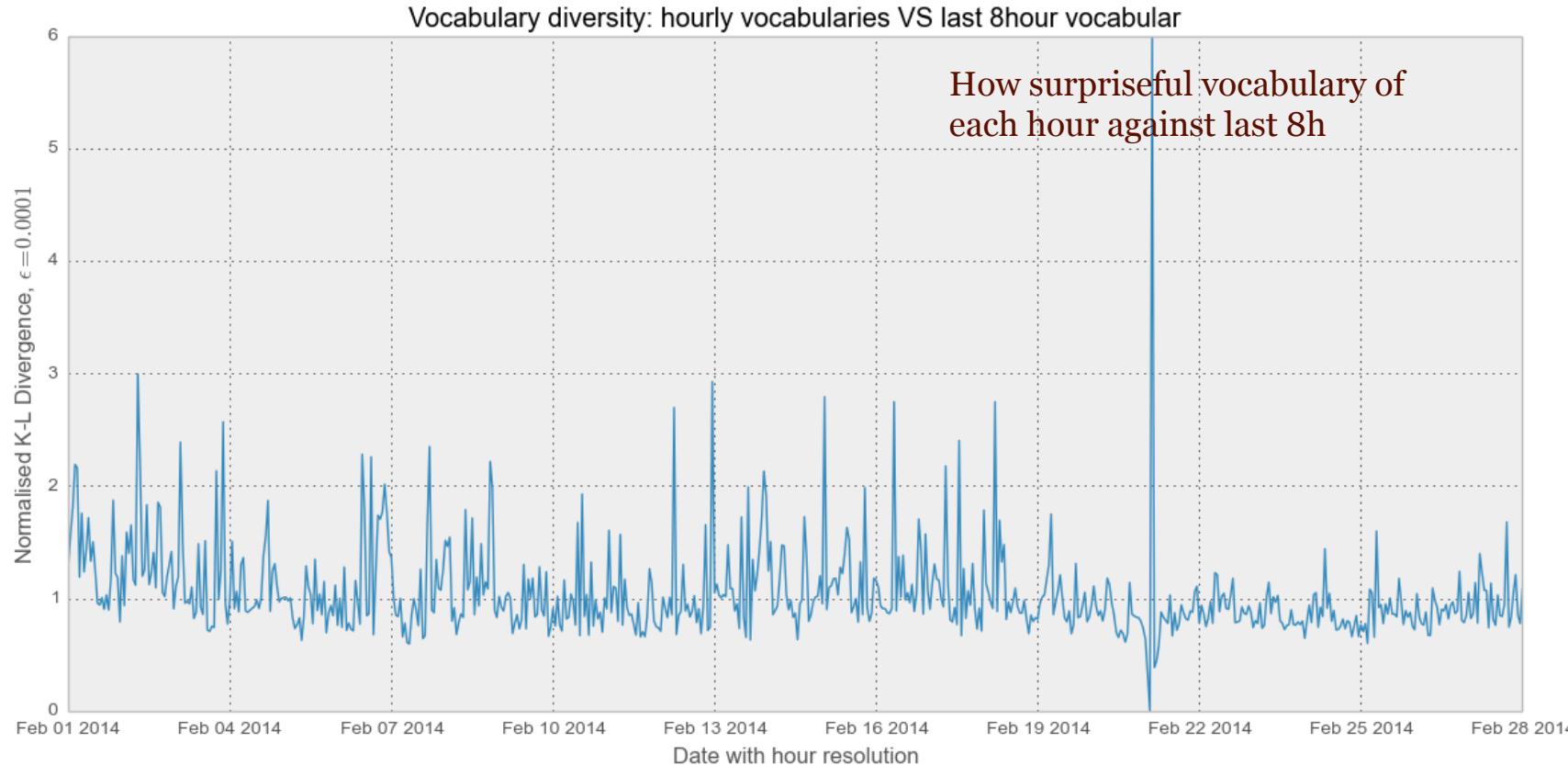
# Vocabulary slots: KLD smoothed



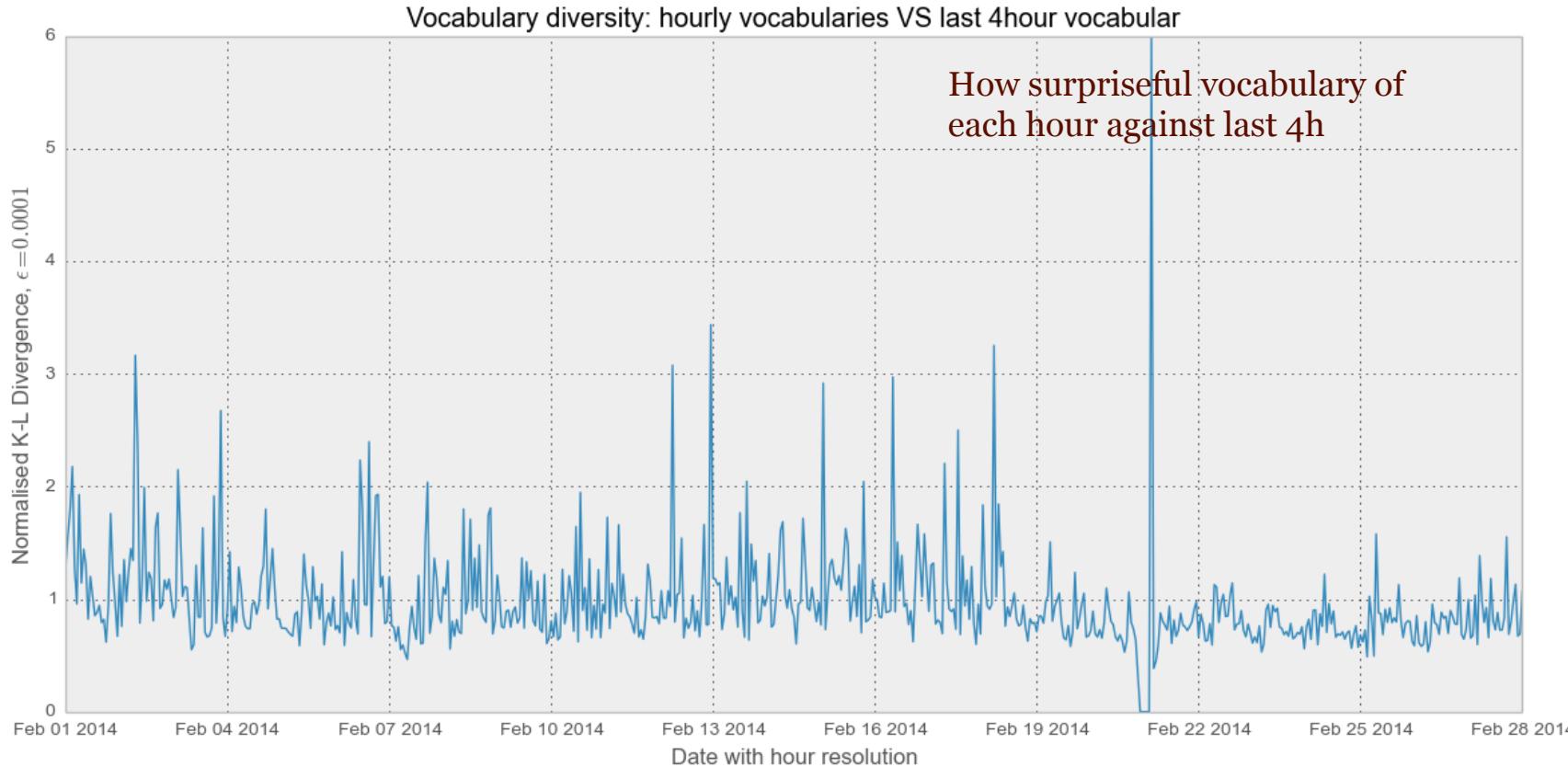
# Vocabulary slots: rolling KLD



# Vocabulary slots: rolling KLD

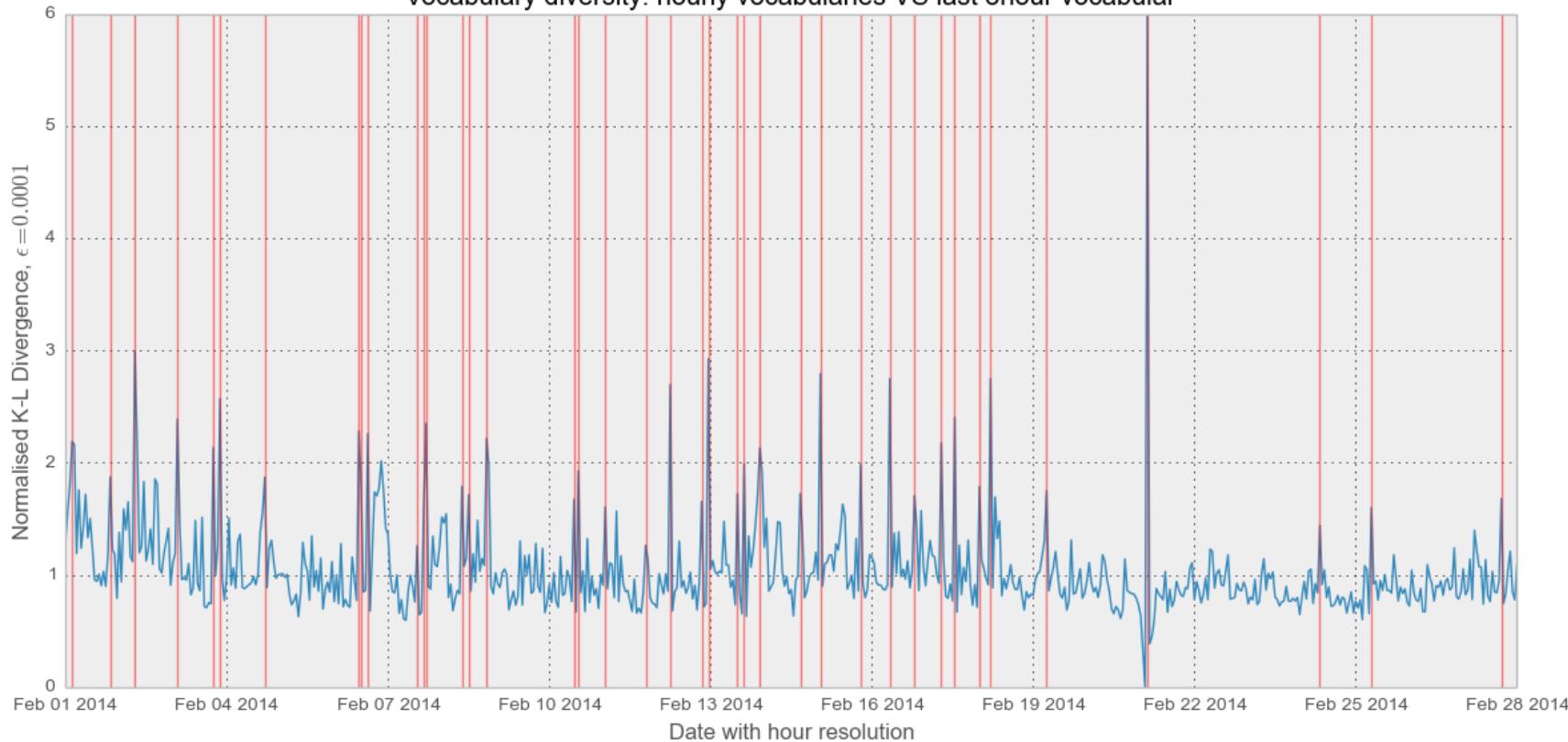


# Vocabulary slots: rolling KLD



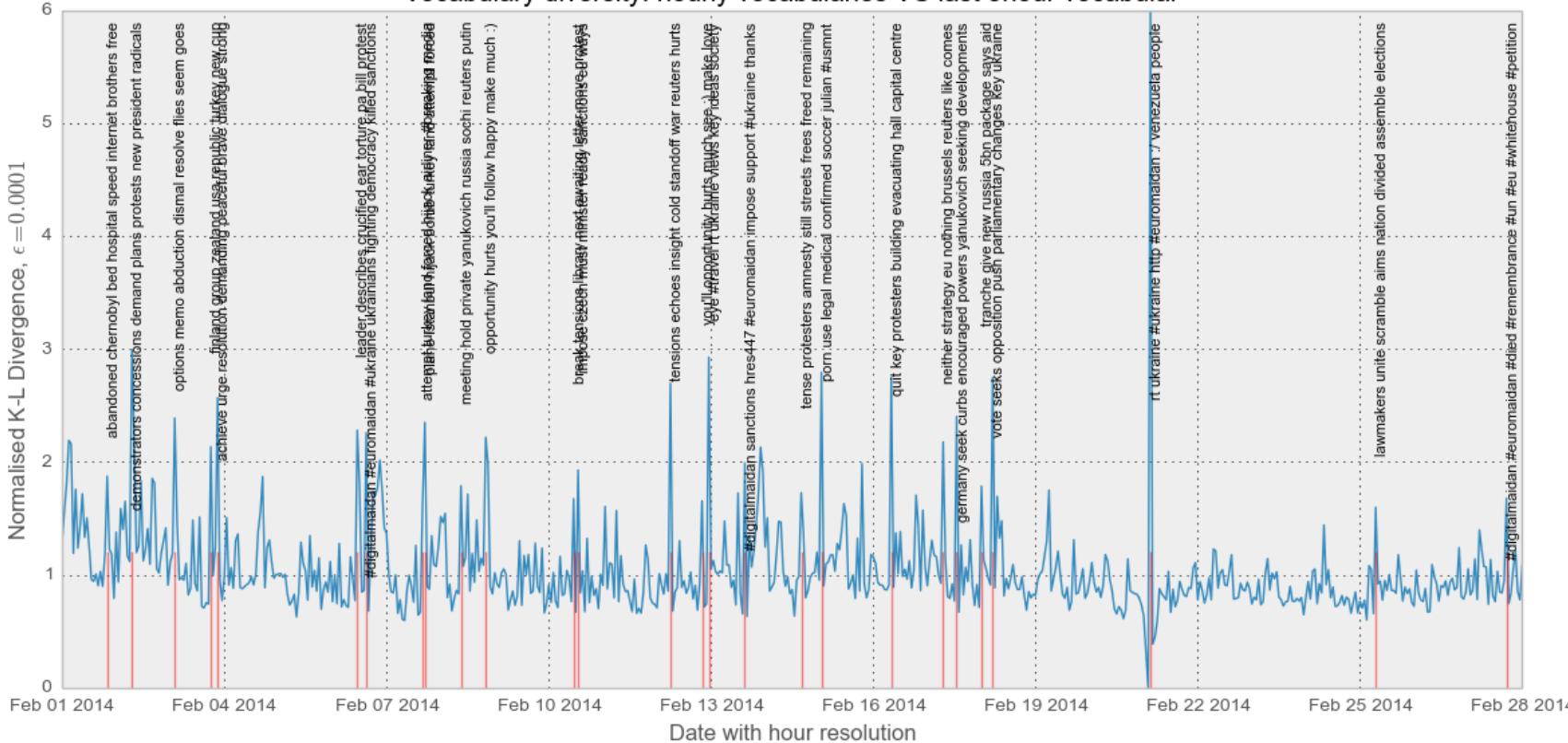
# Rolling KLD outliers

Vocabulary diversity: hourly vocabularies VS last 8hour vocabulary

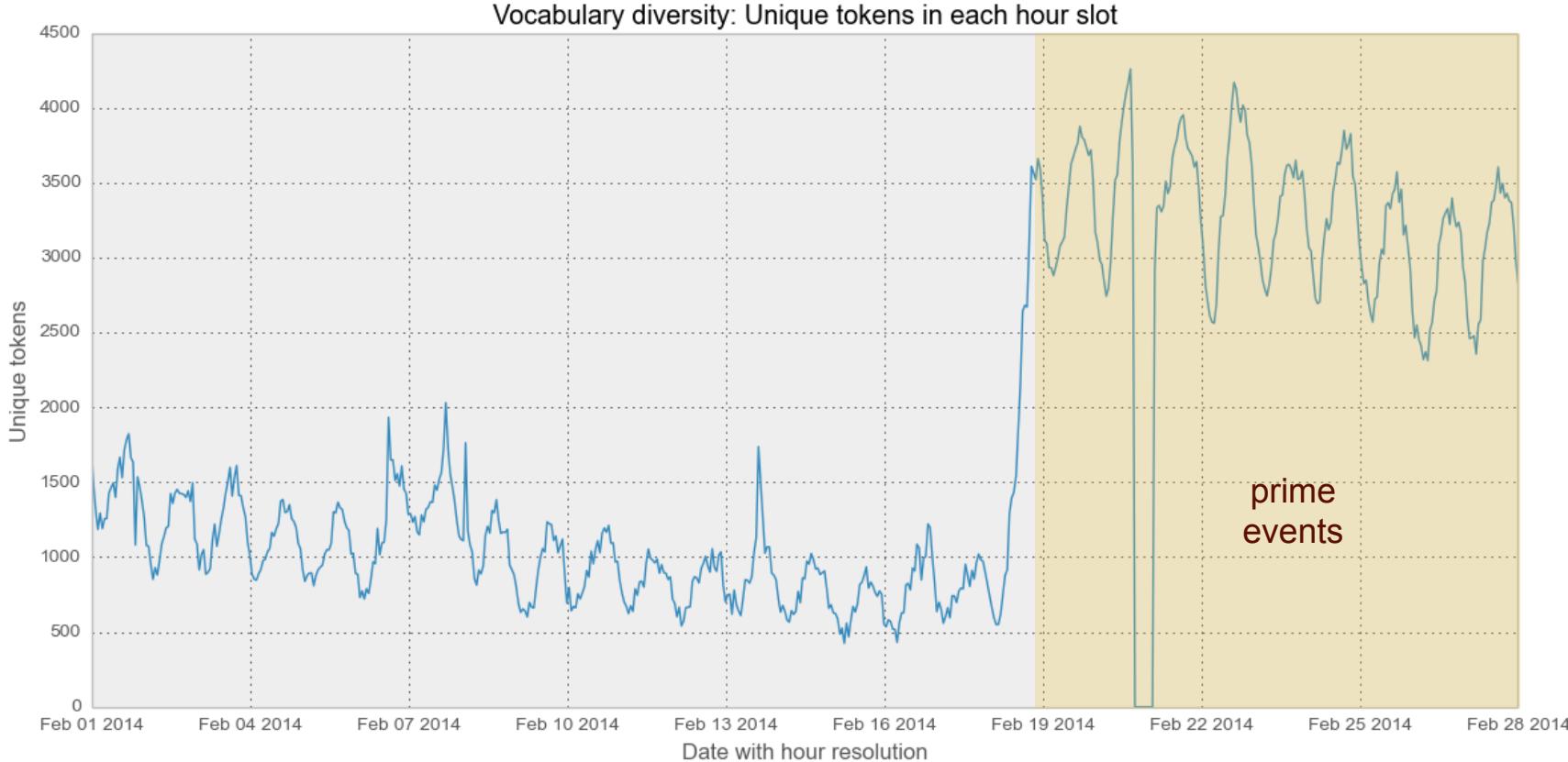


# Rolling KLD outliers tokens

Vocabulary diversity: hourly vocabularies VS last 8hour vocabular

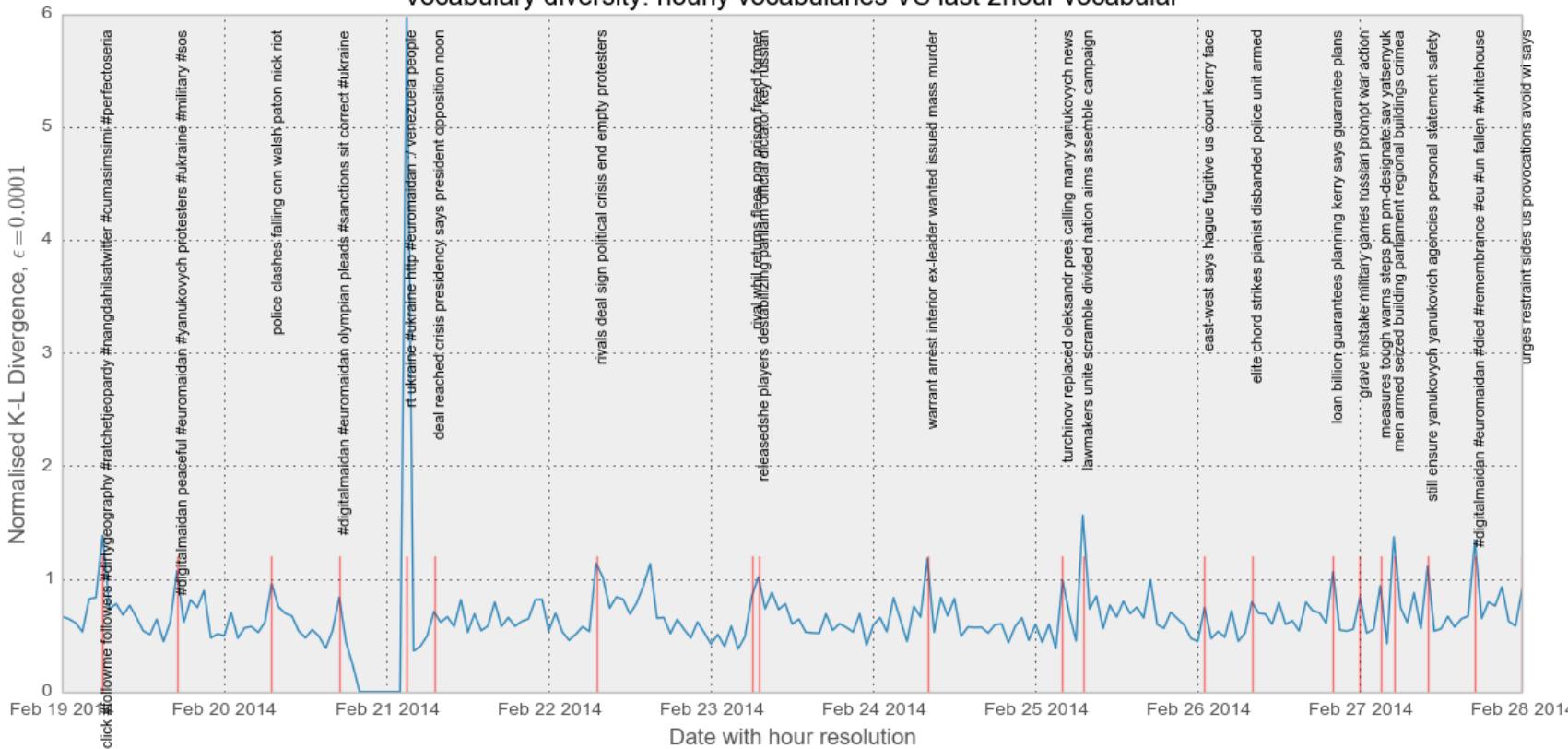


# Further dataset limitation



# Rolling KLD outliers: Feb 19-28

Vocabulary diversity: hourly vocabularies VS last 2hour vocabular



# Surpriseful tweets link ➡

9p.org.ua/tweets x 9p.org.ua/tweets/t1.html



7:00 AM  
February 20, 2014

**Intense clashes. Riot police falling back. Petrol bombs, tear gas, a lot of smoke. #Kiev #Ukraine - via @\_DuncanC http://t.co/WfF3YKR8UD [abiesepr 24:164]**

Intense clashes. Riot police falling back. Petrol bombs, tear gas, a lot of smoke. #Kiev #Ukraine http://t.co/LnKfZVjjs - via @\_DuncanC [olINANNEWS 1499:257]

RT @BBCWorld: Intense clashes. Riot police falling back. Petrol bombs, tear gas, a lot of smoke. #Kiev #Ukraine - via @\_DuncanC http://t.co...  
[@punkobnoi 79:218] →

@atomkigler: I stand with the people of the #Ukraine #PJNET #tcoot #ORPUW #ccot #teaparty #RedNationRising http://t.co/nXRMtcbD# 2A [@Metz\_Byron 749:1026]

This is a special #coffee for #Ukraine #bordia #goodmorning #buenosdias #café http://t.co/ZWP1GBGFG [valeronline 10:49]

"It's an imp. time for Canada to stand for a Ukraine that's democratic, open & transparent." Cdn. Ukrainian Congress http://t.co/aeBpvFlZmY [@Vida\_Jay 1621:2002]

5:00 PM  
February 20, 2014

Thank you  
@CBNews for  
using correct  
spelling of KYIV,  
#EuroMaidan  
#DigitalMaidan  
#Ukraine  
@nytimes  
@nytimesworld  
@politic oABC  
@abcNews  
[gMiller27  
70:60]

RT @RT\_com: MORE: #Ukraine's Yanukovich wanted for mass murder of peaceful citizens - interim interior minister http://t.co/Zr5CGAcVln true?? [...] @akatulkhajuria 103:6551

Russian official: West  
destabilizing Ukraine:  
Ex-PM Yushchenko

CLICK 500  
@Marcorubio #SO  
#Ukraine #Yanuk  
#murders peacef

Intense clashes. Riot  
police falling back. Petrol  
bombs, tear gas, a lot of

BREAKING: #Ukraine  
presidency website says  
deal reached with EU.

Ukraine Rivals Sign Deal  
To End Political Crisis -  
Ukraine Has Deal

Ukraine President flees  
Kiev; rival freed from  
prison. Former PM Yulia

FEB. 20 FEB. 21 FEB. 22 FEB. 23 FEB. 24

# Surpriseful tweets link ➔

9p.org.ua/tweets x  
9p.org.ua/tweets/t2.html

3:00 AM  
RT  
@jonathanhainin:  
Masha Lipman  
on Ukraine's  
"brittle  
nationhood" and  
"the related risk  
of disintegration":  
<http://t.co/HgUc1Bwz6y> [6412:1183] [retweet](#)

BREAKING: Ukraine says deal reached with EU, Russia on settlement of crisis <http://t.co/QNn5s2U3h1>  
<http://t.co/Z58AN6kv7P> (@globalnews 3256:1780)

7:00 AM  
February 21, 2014  
**BREAKING: #Ukraine presidency website says deal reached with EU, Russia on settlement of political crisis.** [[@CFRAOttawa 14255:96](#)]

Breaking News: Ukraine presidency website says deal reached with EU, Russia on settlement of political crisis. [[@WOKVNews 6074:437](#)]  
Breaking (AP) -- Ukraine presidency website says deal reached with EU, Russia on settlement of political crisis. [[@FOX29philly 85178:3496](#)]

This is fantastic. \*@anneapplebaum: My Ukrainian lexicon: clichés to treat skeptically when speaking of Ukraine <http://t.co/yHq3Cjqlq>\* [[@pettybooshwh 1688:169](#)]  
Wassap Ibadan  
Trends#TIMINUkraineVenny#EveryMansDeliverance#HappilyEverAfterBy\_Roland#KongaFashionSanusiGEJNigeria#RepYourHood [[@WASSAP\\_IBADAN 3130:651](#)]  
@TelegraphNews: Financial crisis threatens Russia as #Ukraine spins out of control | <http://t.co/jiVC5oHKSQ>\*  
@AmbroseEP on geopolitical woe [[@domcavendish 6768:1357](#)]

CLICK 500  
@SenBobCorker  
#SOS #Ukraine  
#Yanukovych

Intense clashes. Riot police falling back. Petrol bombs, tear gas.

Thank you @CBC for using correct spelling of KYIV.

RT  
@jonathanhainin:  
Masha Lipman on

Ukraine Rivals Sign Deal To End Political Crisis - Ukraine Has

Ukraine President flees Kiev; rival freed from prison: Former Ex-PM

#Ukraine arrest warrant for ex-leader

FEB. 20 FEB. 21 FEB. 22 FEB. 23 FEB. 24 FEB. 25

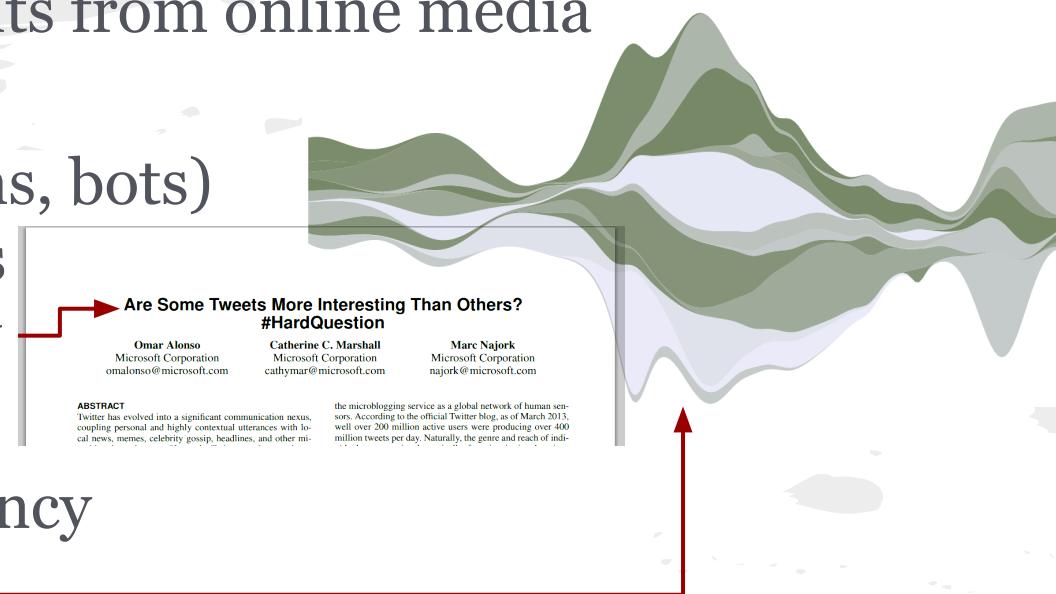
Only from users with  
+500followers



[majdannezalezhnosti.blogspot.com](http://majdannezalezhnosti.blogspot.com)

# To improve

1. Benchmark
  - a. export 'hot' events from online media
2. Fight bots
  - a. spam (repetitions, bots)
  - b. 'forced' opinions
  - c. filter low quality
3. Topic model
  - a. no Term Frequency
  - b. split topics (!)



Q?