

# Mining Social Media

Carlos Castillo

European Summer School in Information Retrieval  
July 2019

## Sources:

Some slides from “Twitter and the Real World” CIKM’13 Tutorial.

See also “Twitter: A Digital Socioscope” (2015) by Mejova, Weber, and Macy  
Ghani et al. “Social media big data analytics: A survey” (2018).

# Douglas Adams on new technologies

Anything that is in the world when you are born is normal and ordinary and just a natural part of the way the world works.

Anything that is invented between when you are 15 and 35 is new and exciting and revolutionary and you can probably get a career in it.

Anything invented after you are 35 is against the natural order of things.

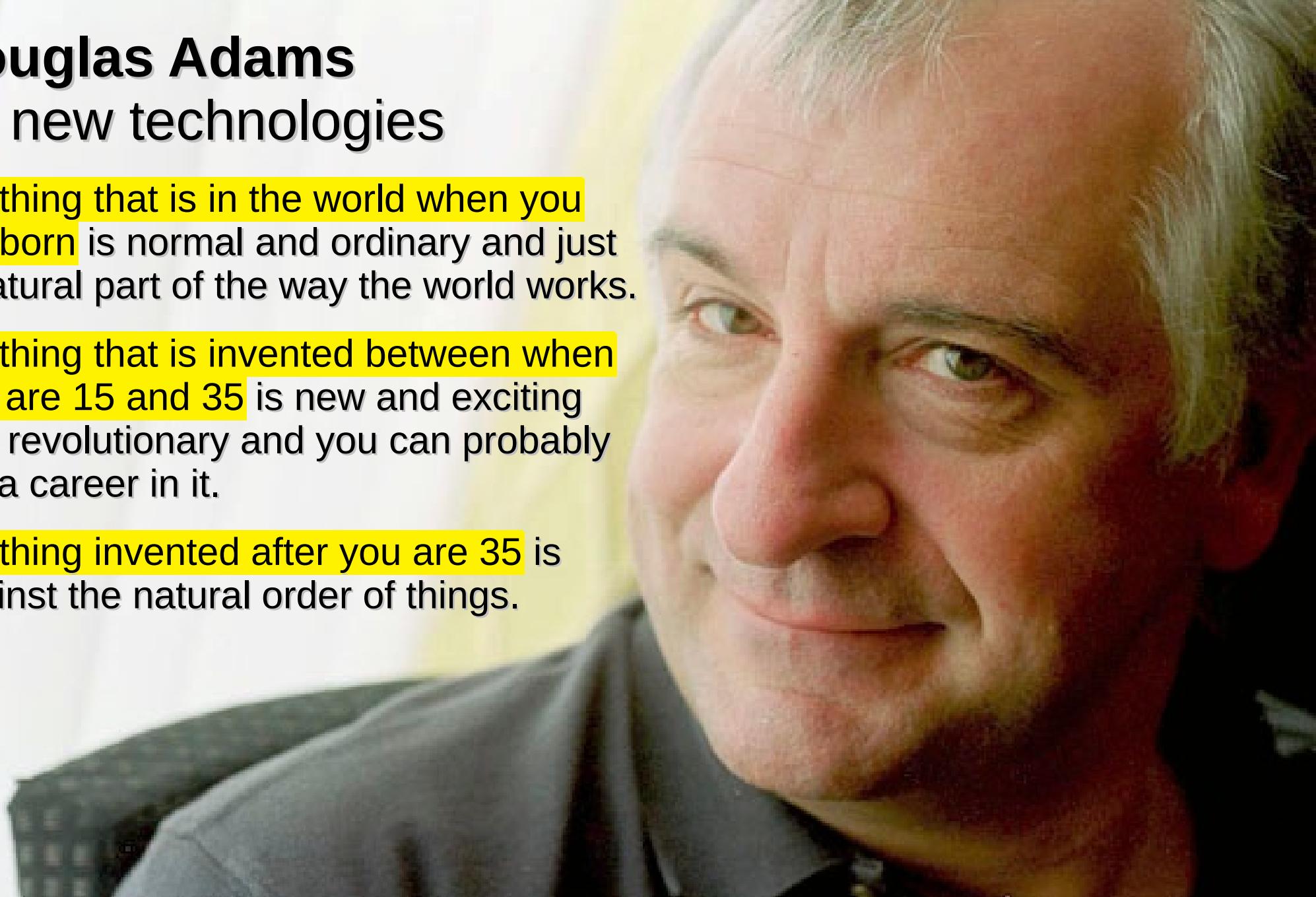


Photo by Michael Hughes - Flickr, CC BY-SA 2.0,  
<https://commons.wikimedia.org/w/index.php?curid=6661797>

Definitions  
Motivations  
Caveats

# Definitions

- **Social software**
  - Software to facilitate or mediate social interactions
- **Social networking sites**
  - Sites allowing (semi-)public profiles, connections, and viewing/interacting with other people's activities  
[\[boyd & Ellison 2007\]](#)
- **User-generated content**
  - Content created by end-users of social software outside of their professional routine

# Why mining social media?

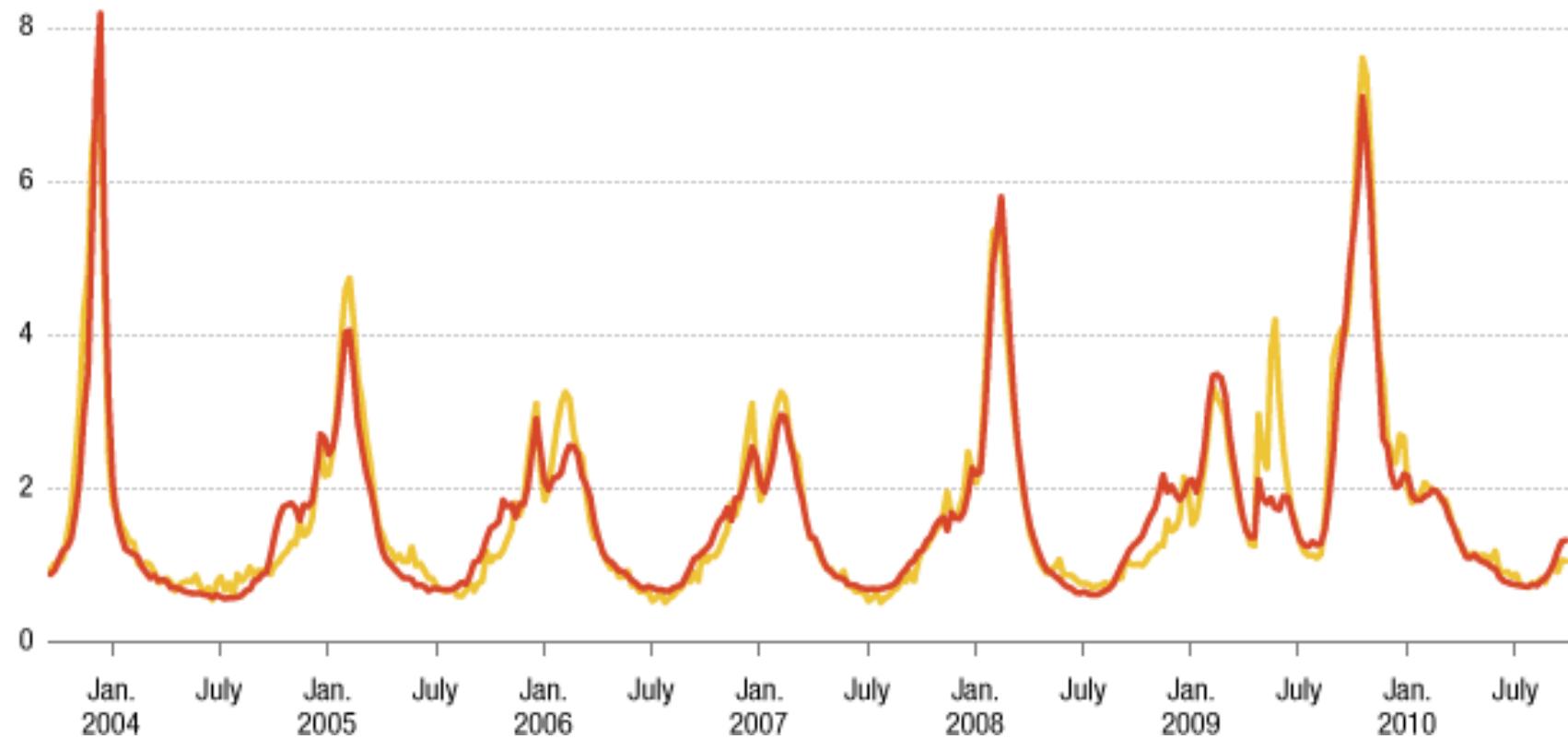
- **Seen as an alternative to traditional opinion polls**
- “*What do people think about X?*”
- “*How do they feel about X?*”
  - First steps: sentiment analysis and opinion mining
- Attractive for many reasons including:
  - Lower latency (waiting time)
  - Lower cost
  - Larger population



# Key example: Google Flu Trends

— Google Flu Tracker    — Official report

10%



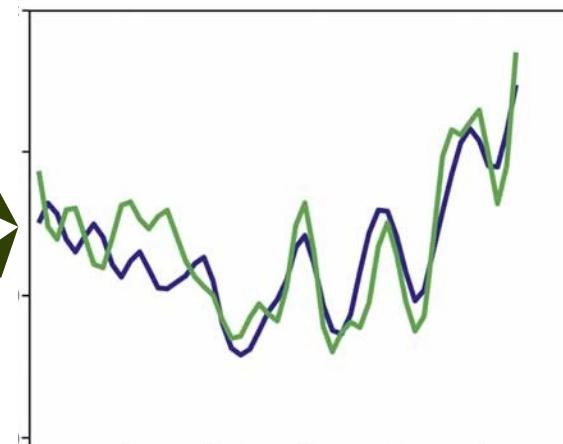
Low latency: people query symptoms **before** going to the doctor!

# Many social media mining papers

## 1. Domain-specific data



### 3. Correlation/Influence



# Profit?

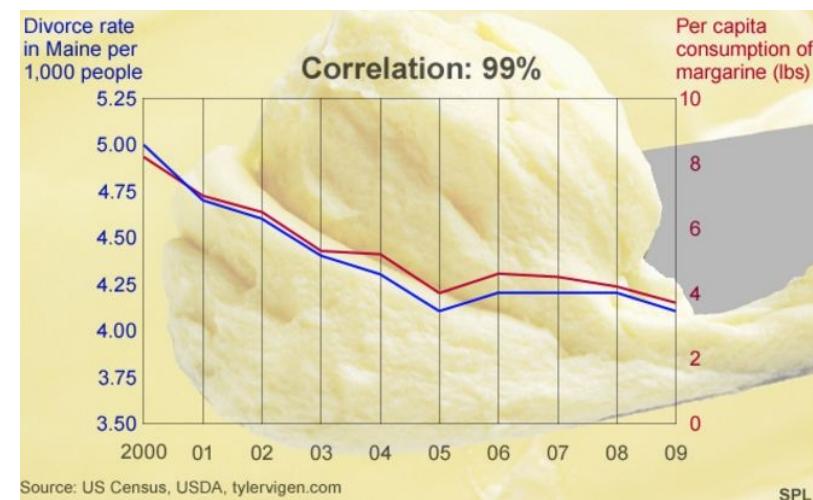
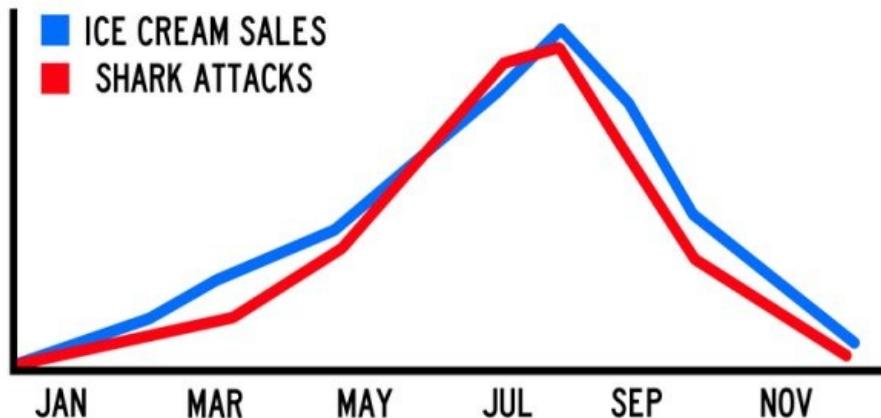
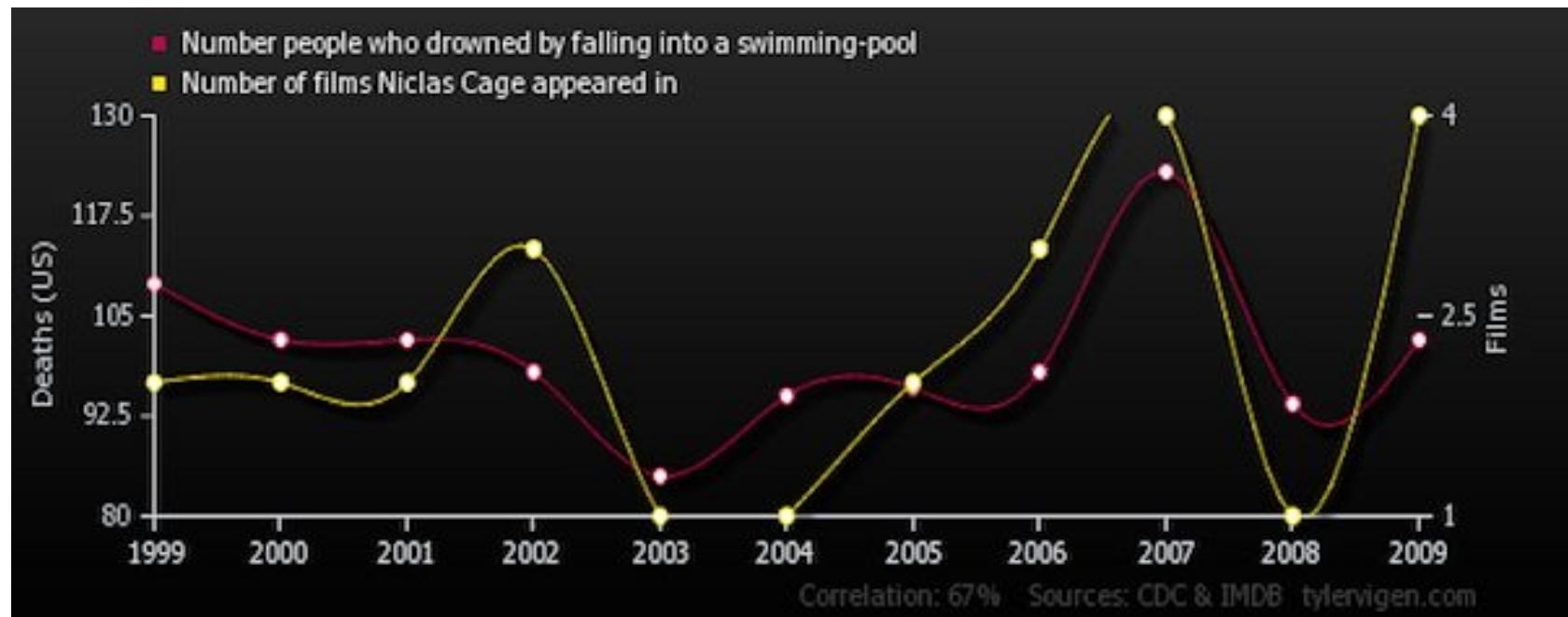
## 2. Social media data



# The devil is in the details

- Which domain-specific data?
  - This is not always readily available
  - Many biases can be introduced (see talk by Ricardo B.-Y.)
- Mapping social media data to a time series?
  - Geolocation of messages
  - Mapping to topics/sentiments/intents or other characteristics
  - What is the variable: Volume? Sentiment? Other?
- Measuring correlation/influence
  - Correlation (lagged); Transfer entropy
- Finding a causal mechanism

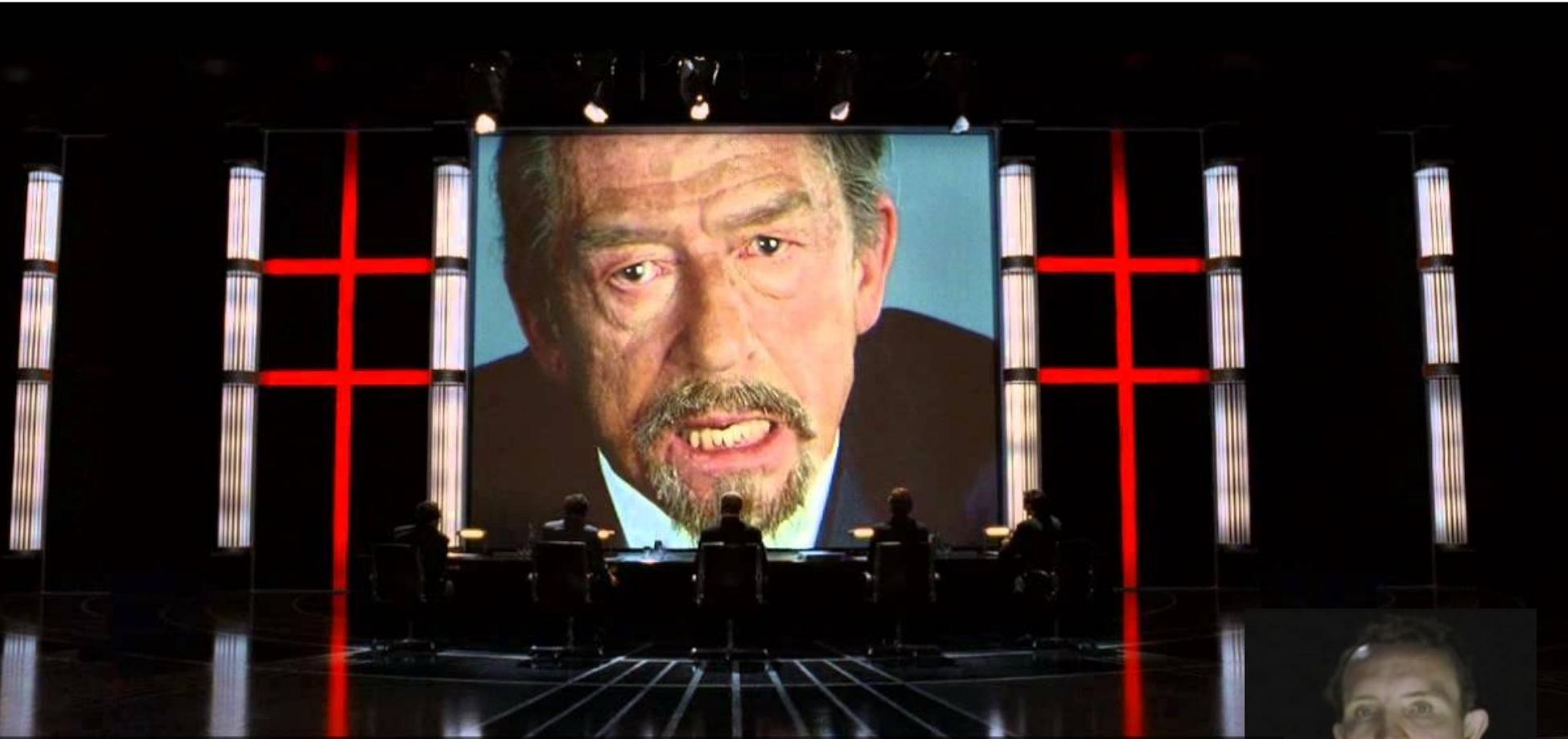
# Caveat 1: correlation might be spurious



# Caveat 2: correlation might be useless

- Sometimes there are much better predictors
- Social media can be used to predict box office revenue
  - But ticket sales on first weekend *almost* always determine total sales, with exceptions  
Citizen Kane (1941), Blade Runner (1982), Fight Club (1999)
- Social media can be used to detect earthquakes/floods
  - But sensors are quite dense in wealthy regions of the world

# Caveat 3: the war on terror and “pink lists” (secondary use)



We also are currently monitoring a lot of phone surveillance indicating a high percentage of conversation concerned with the explosions.

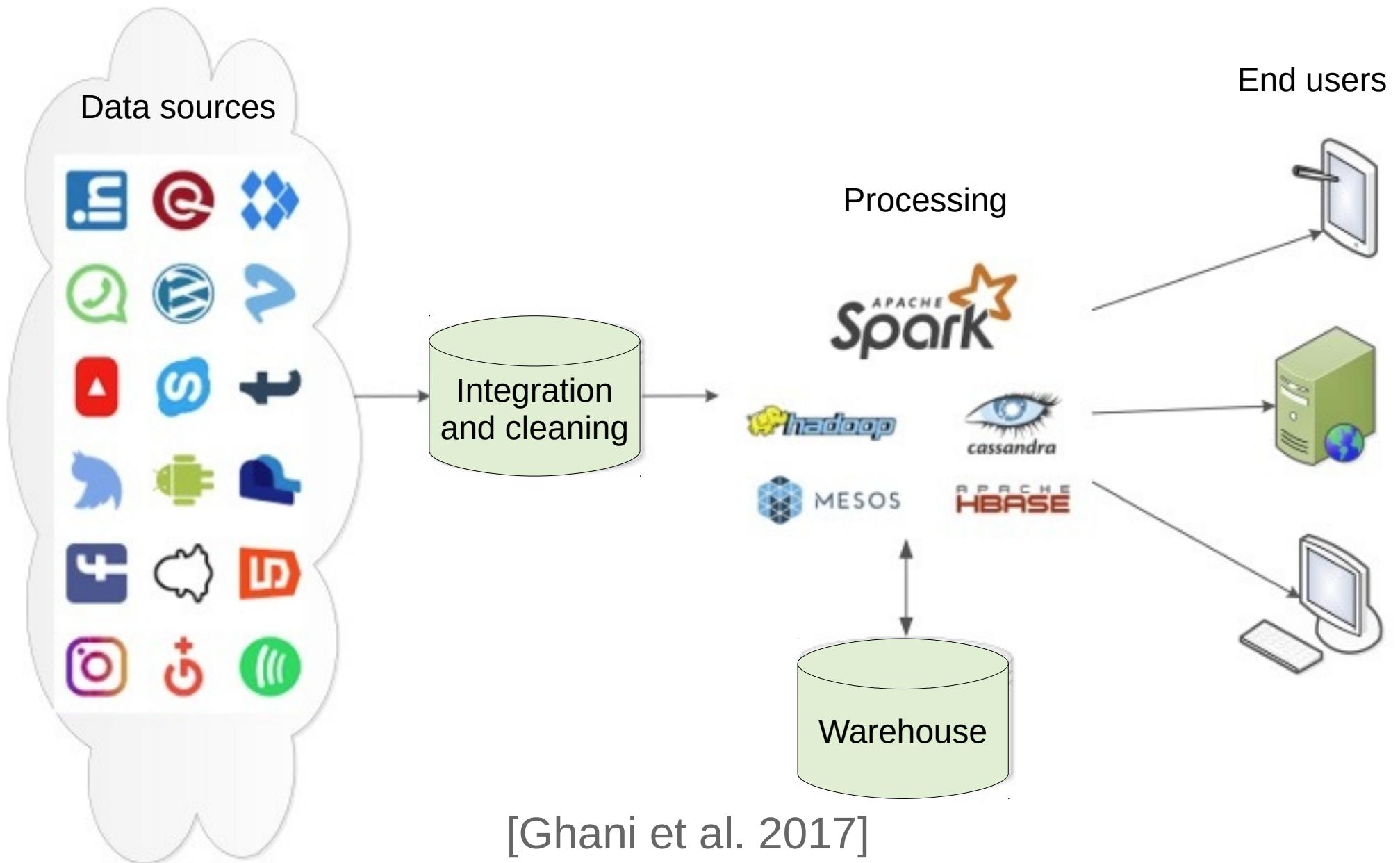


# Applications Methods Architectures

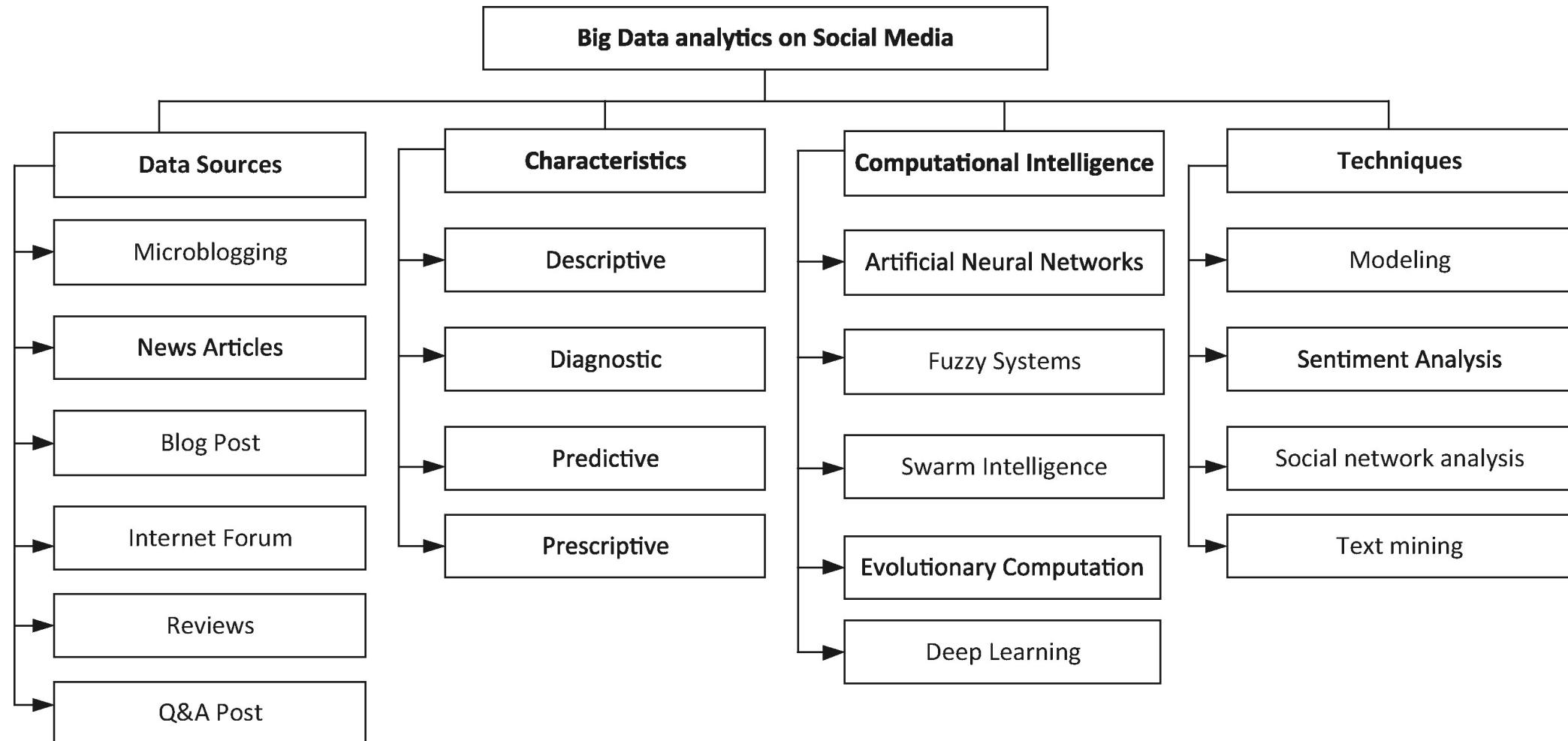
# Analysis of social media

- Most common **applications**
  - trend discovery, social media analytics, sentiment analysis, and opinion mining
- Most common **methods**
  - machine learning, particularly supervised classification and deep learning methods
- Most common **architectures** (see next slide)

# Architectures for SM analysis

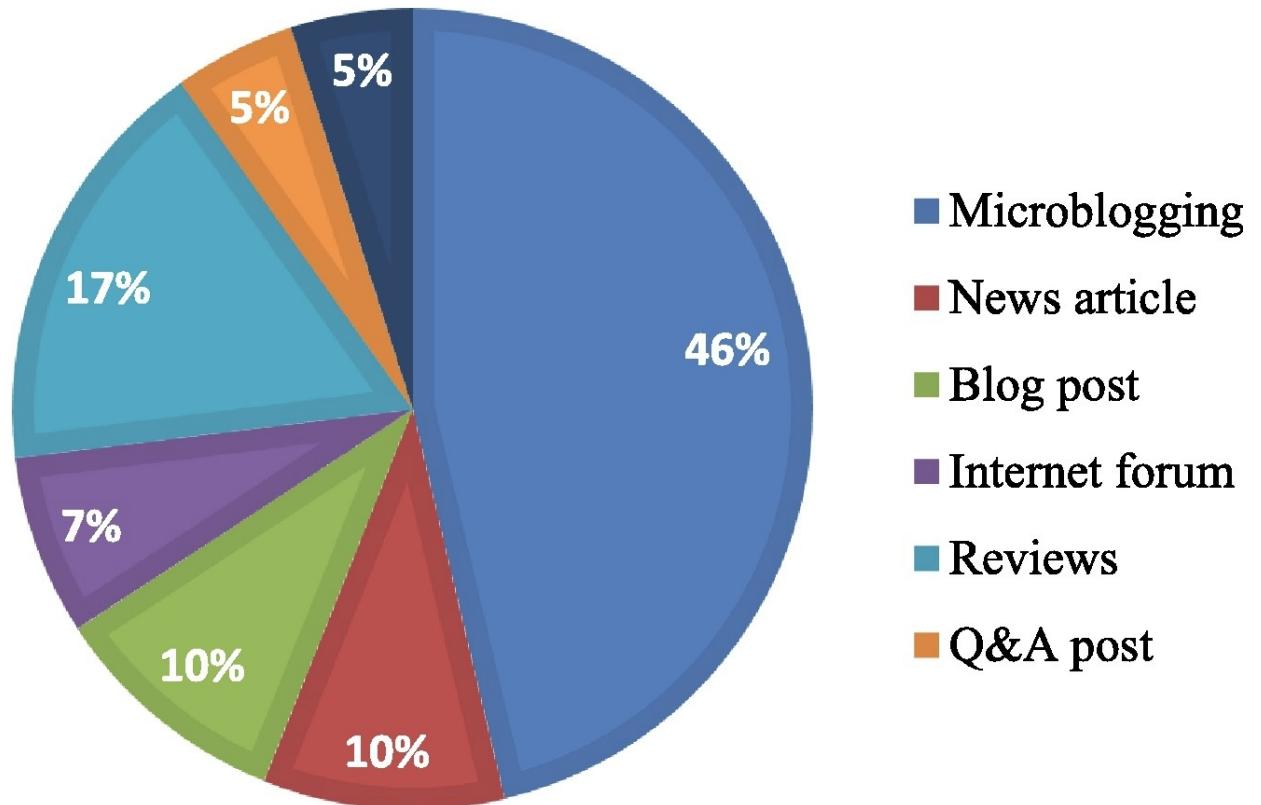


# Elements of SM analytics

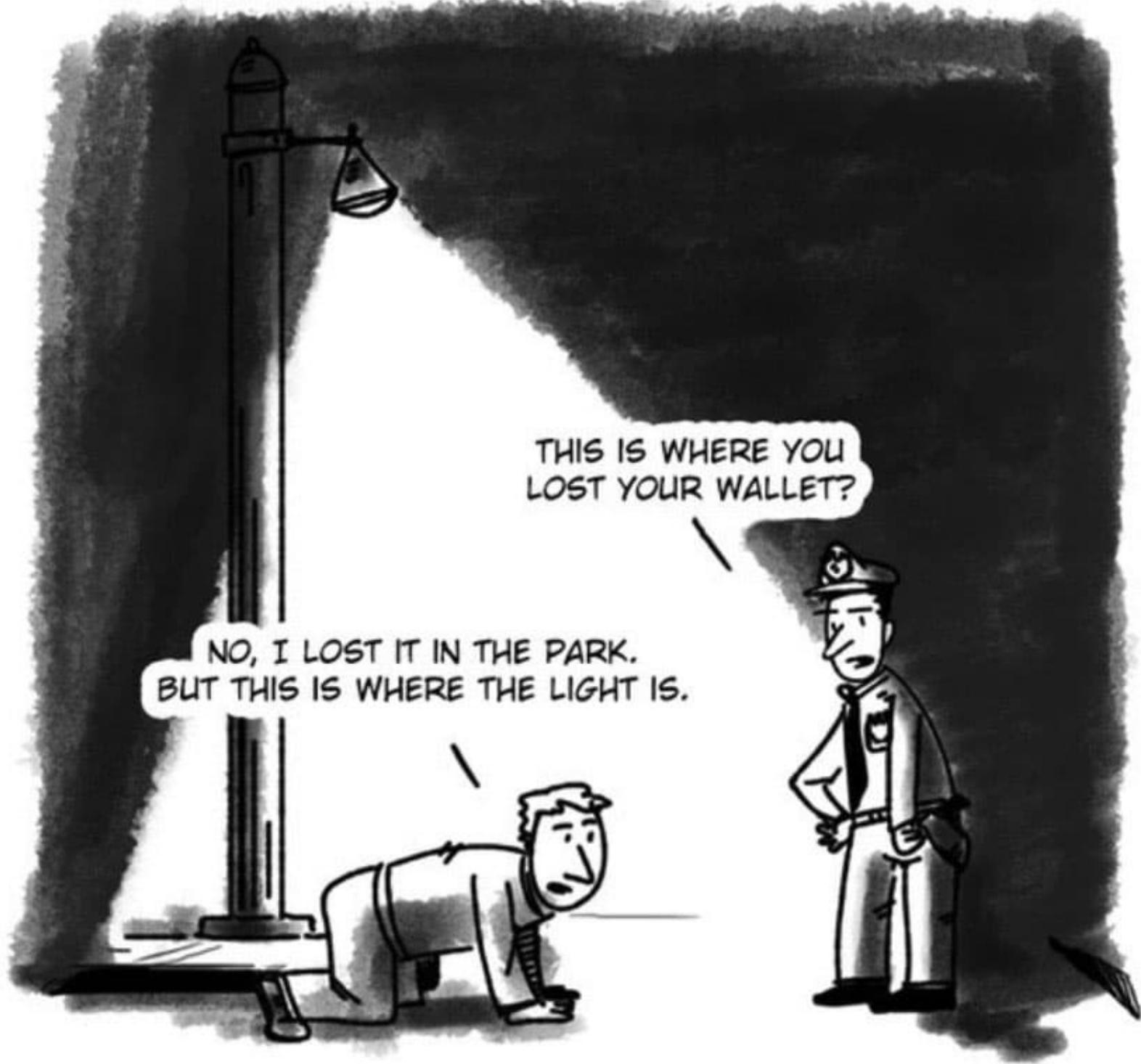


# Data sources

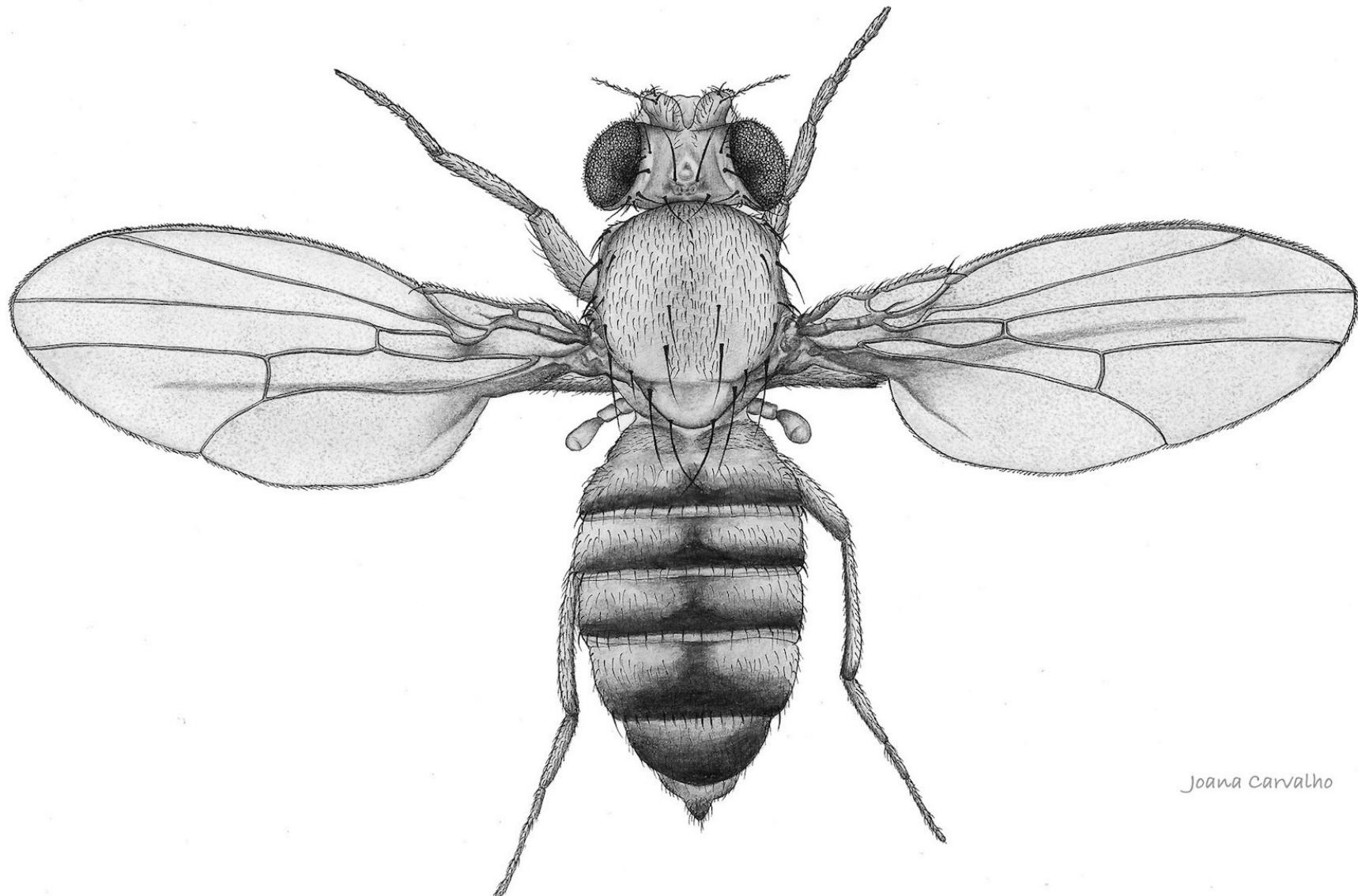
Big Data analytics on Social Media



[Ghani et al. 2017]



Twitter = *Drosophila melanogaster*



# Types of SM analysis

- **Descriptive**
  - Observe and report on **what** happens
- **Diagnostic**
  - Explain **why** something happens
- **Predictive**
  - Predict what **will** continue happening or happen
- **Prescriptive**
  - Indicate what **should** happen to achieve an outcome

# Common techniques

- **Modeling user behavior**
  - capturing, understanding, describing, ...
- **Sentiment analysis**
  - positive, neutral, negative, ...
- **Social network analysis**
  - clusters, centrality, anomaly detection, ...
- **Text mining methods**
  - extraction, summarization, classification, ...
- **Natural experiment methods / observational studies**
  - difference-in-differences, regression discontinuity, matching, ...

# Sentiment analysis: Valence and Arousal

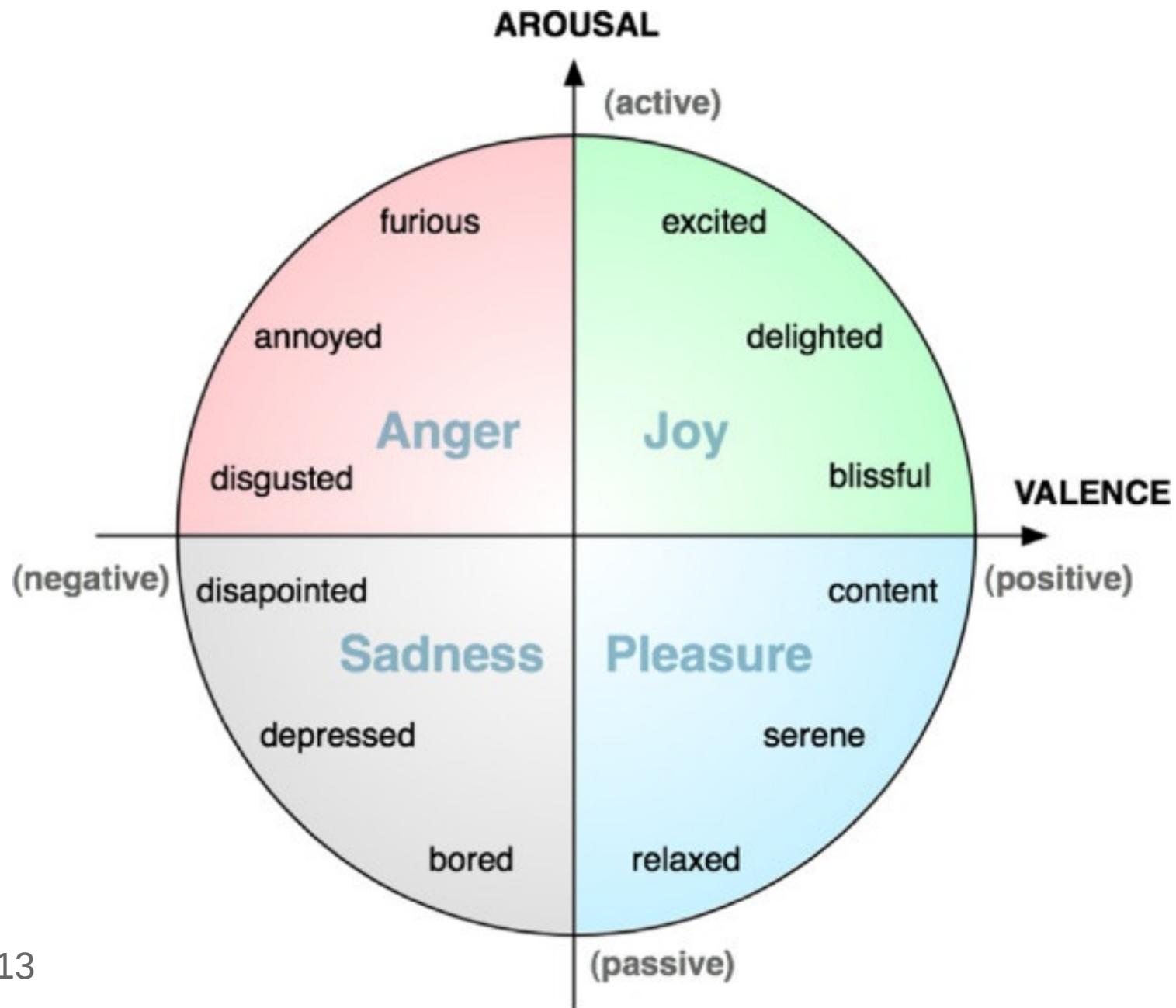
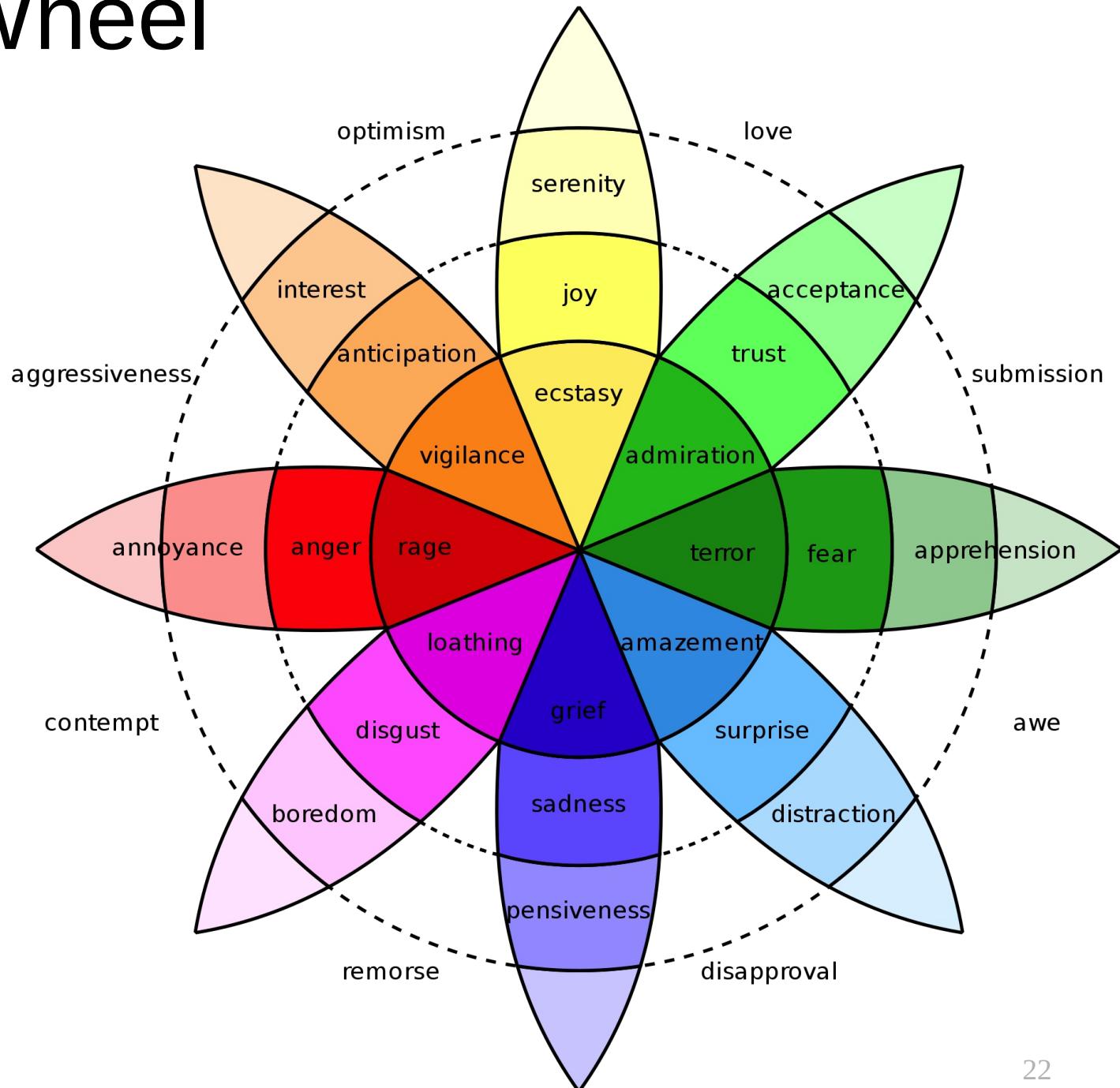


Figure: Yazdani et al. 2013

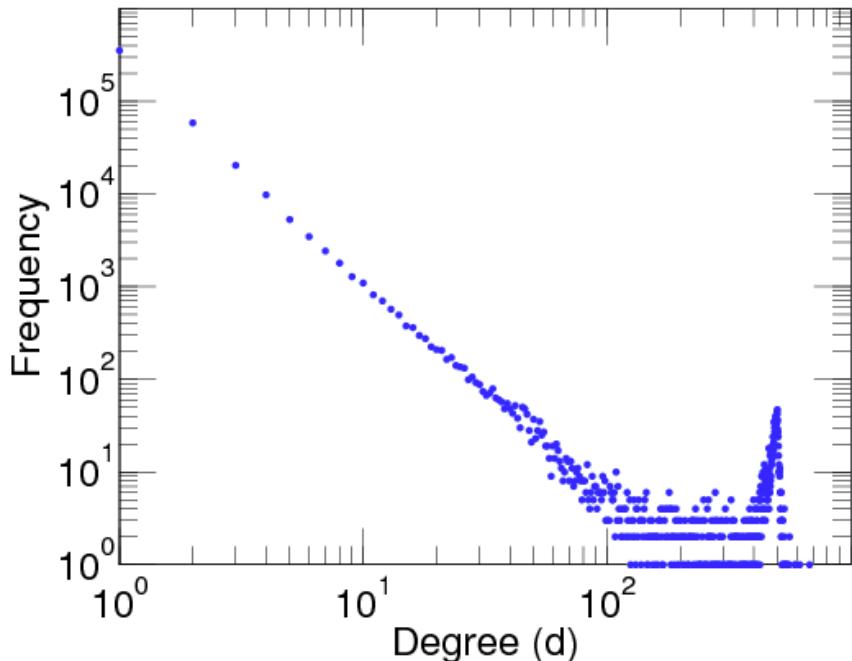
# Sentiment analysis: Plutchik's Wheel



# Network analysis 1/2

## social networks are scale-free networks

Twitter [De Choudhury et al. 2010]

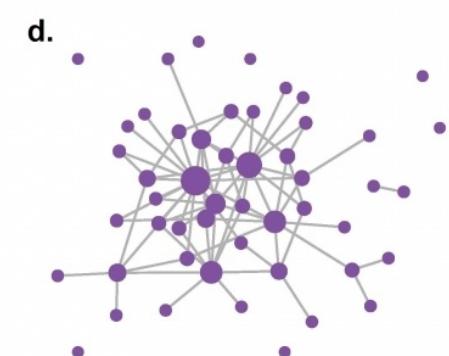
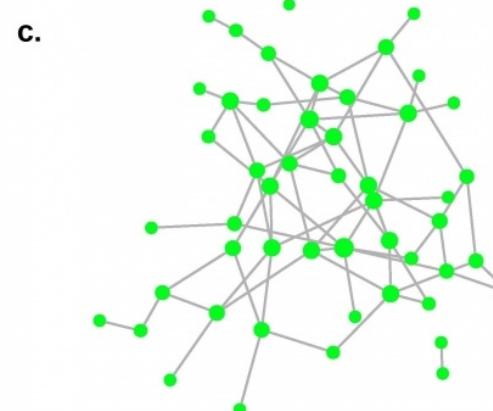
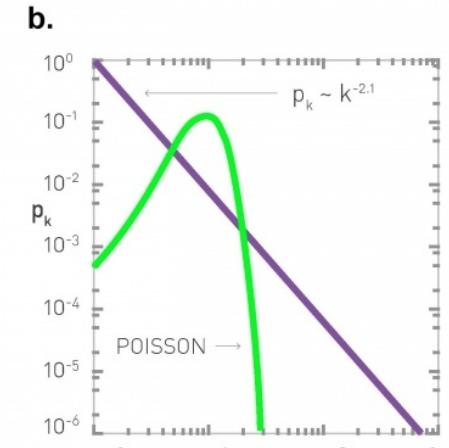
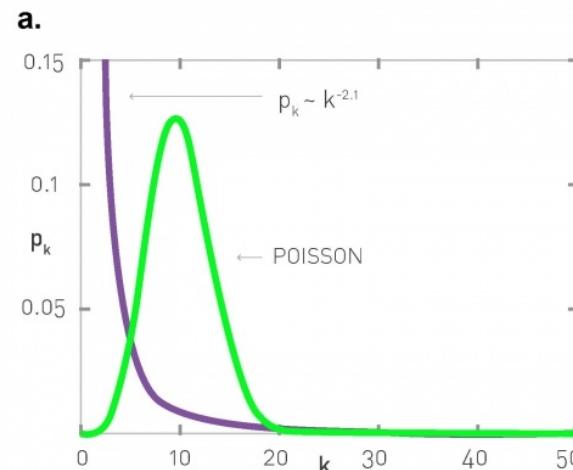


Similar phenomena are observed with other measures of *centrality*

Random (Erdős–Rényi model)

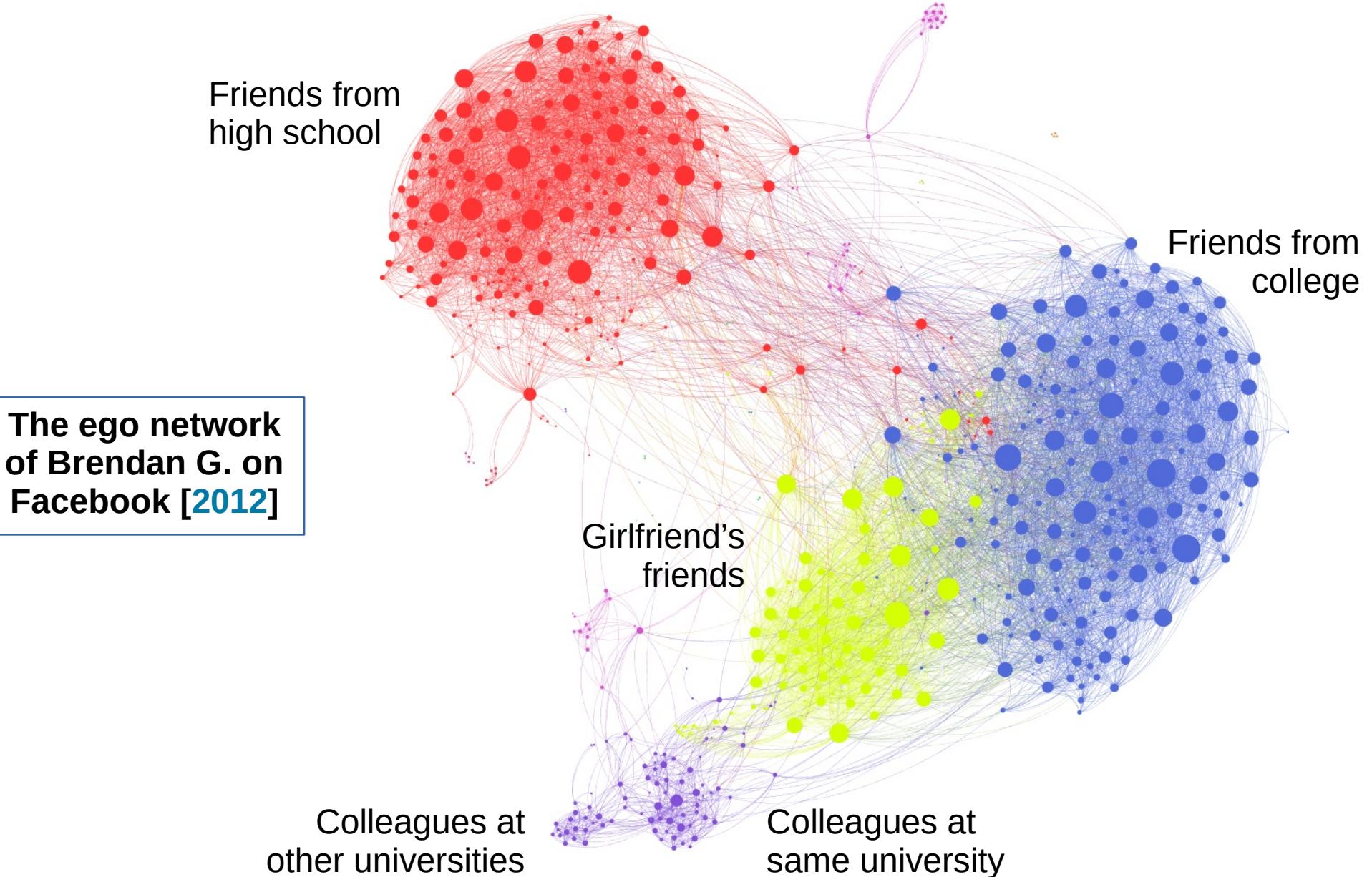
VS.

Scale-free (Barabási-Albert model)



# Network analysis 2/2

social networks have **community structure**



# Examples

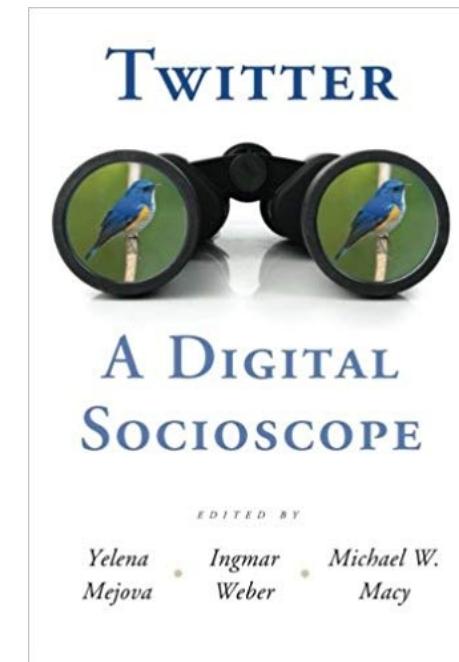
# ICWSM-2019



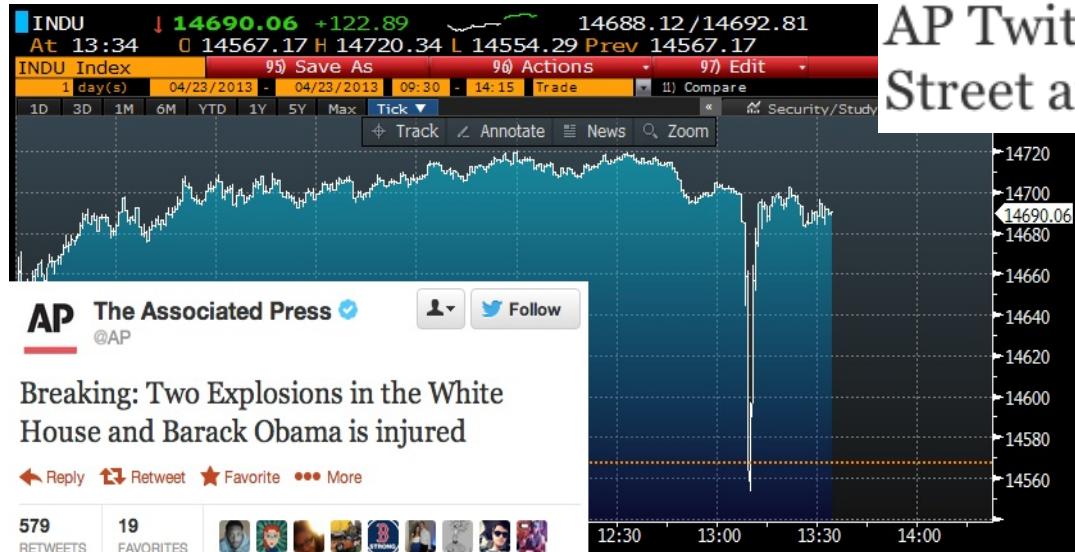
# Tag cloud of titles of accepted papers

# Typical social media mining topics

- **Economics:** marketing, stock market, ...
- **Politics:** elections, demonstrations, ...
- **Public health:** epidemics, ...
- **Smart cities:** transportation, pollution, ...
- **Disasters and mass convergence events**



# When Twitter sneezes, the stock market gets the flu ...



AP Twitter hack causes panic on Wall Street and sends Dow plunging

*During those three minutes, the "fake tweet erased \$136 billion in equity market value,"*  
- Bloomberg News

<http://www.washingtonpost.com/blogs/worldviews/wp/2013/04/23/syrian-hackers-claim-ap-hack-that-tipped-stock-market-by-136-billion-is-it-terrorism/>

## Twitter Death Rumor Leads to Spike in Oil Prices



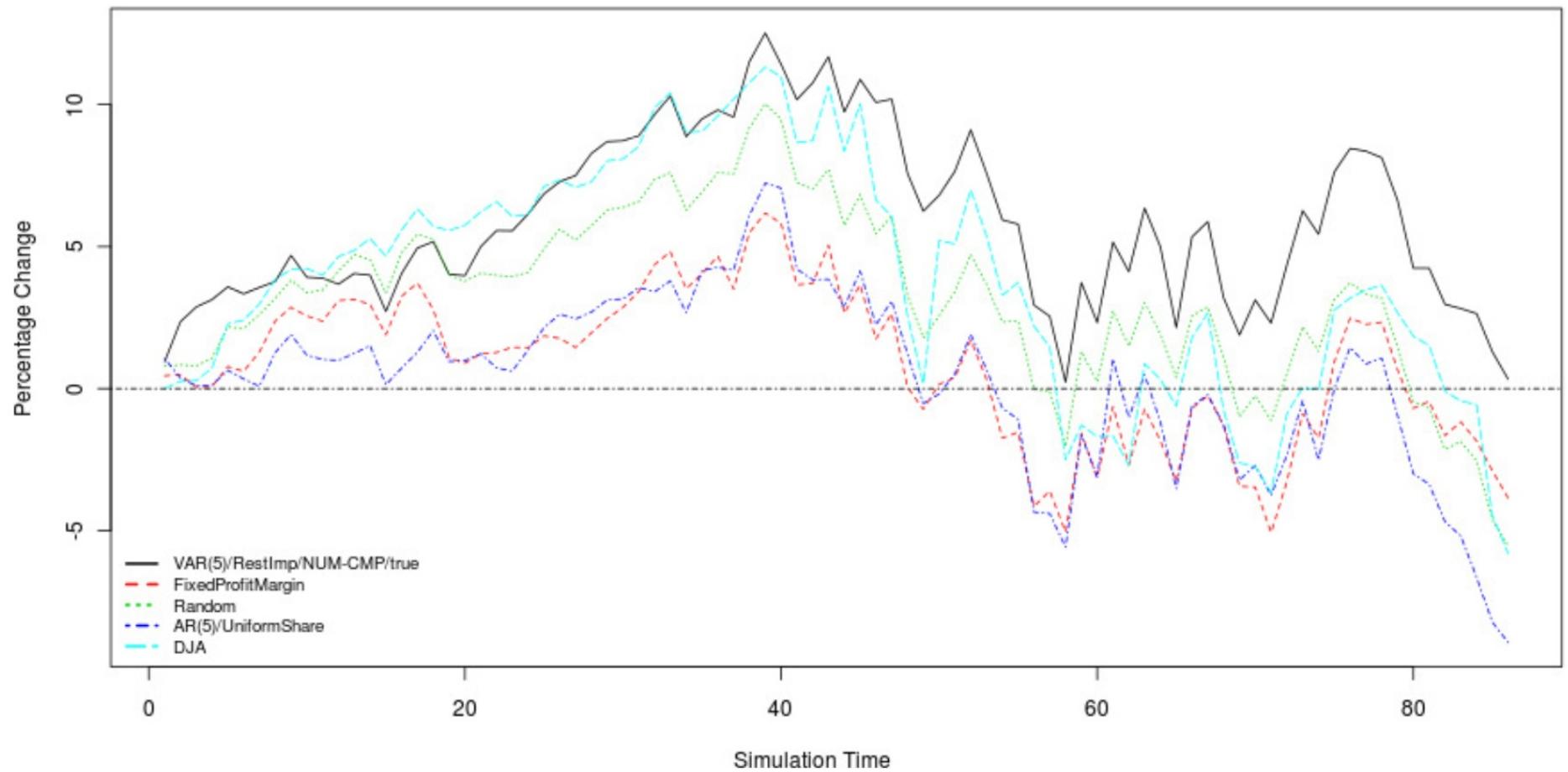
<http://mashable.com/2012/08/07/twitter-rumor-oil-price/>

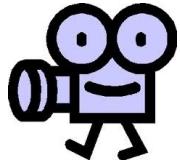
## Netflix CEO's Facebook Post Triggered SEC Wells Notice



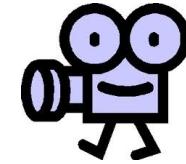
<http://www.cnbc.com/id/100289227>

# Trading stock using social media



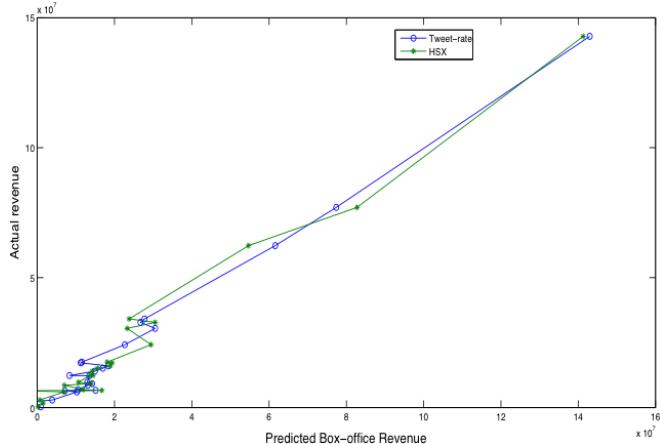


# Movies!

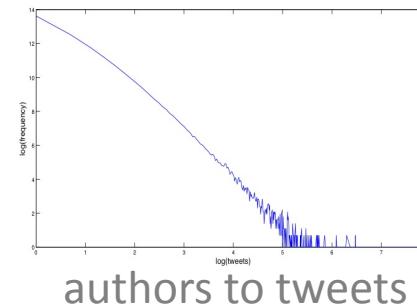


- 2.89 million tweets
- 24 movies (manually compiled keywords)

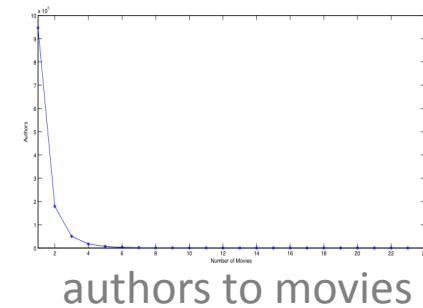
**Correlation** = 0.90  
 (tweet volume vs. box office gross revenue)



predicted vs actual box office scores



authors to tweets



authors to movies

linear regression using previous week's tweets to predict weekend box office gross:

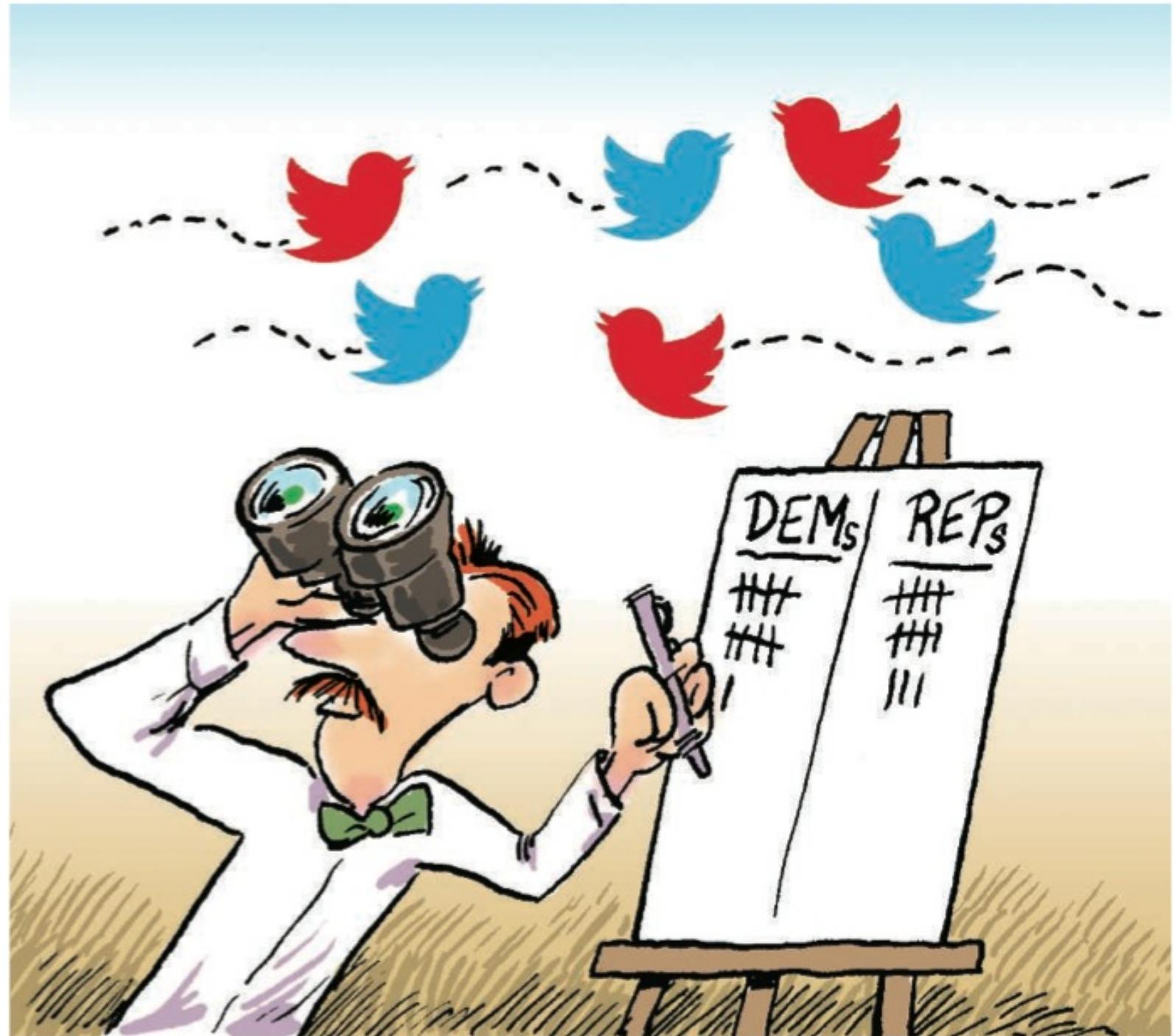
	Adj R <sup>2</sup>
Average Tweet-rate	0.80
Tweet-rate time series	0.93
<b>Tweet-rate time series + theater count</b>	<b>0.973</b>
HSX time series + theater count	0.965



DEMOCRAT  
≈ left



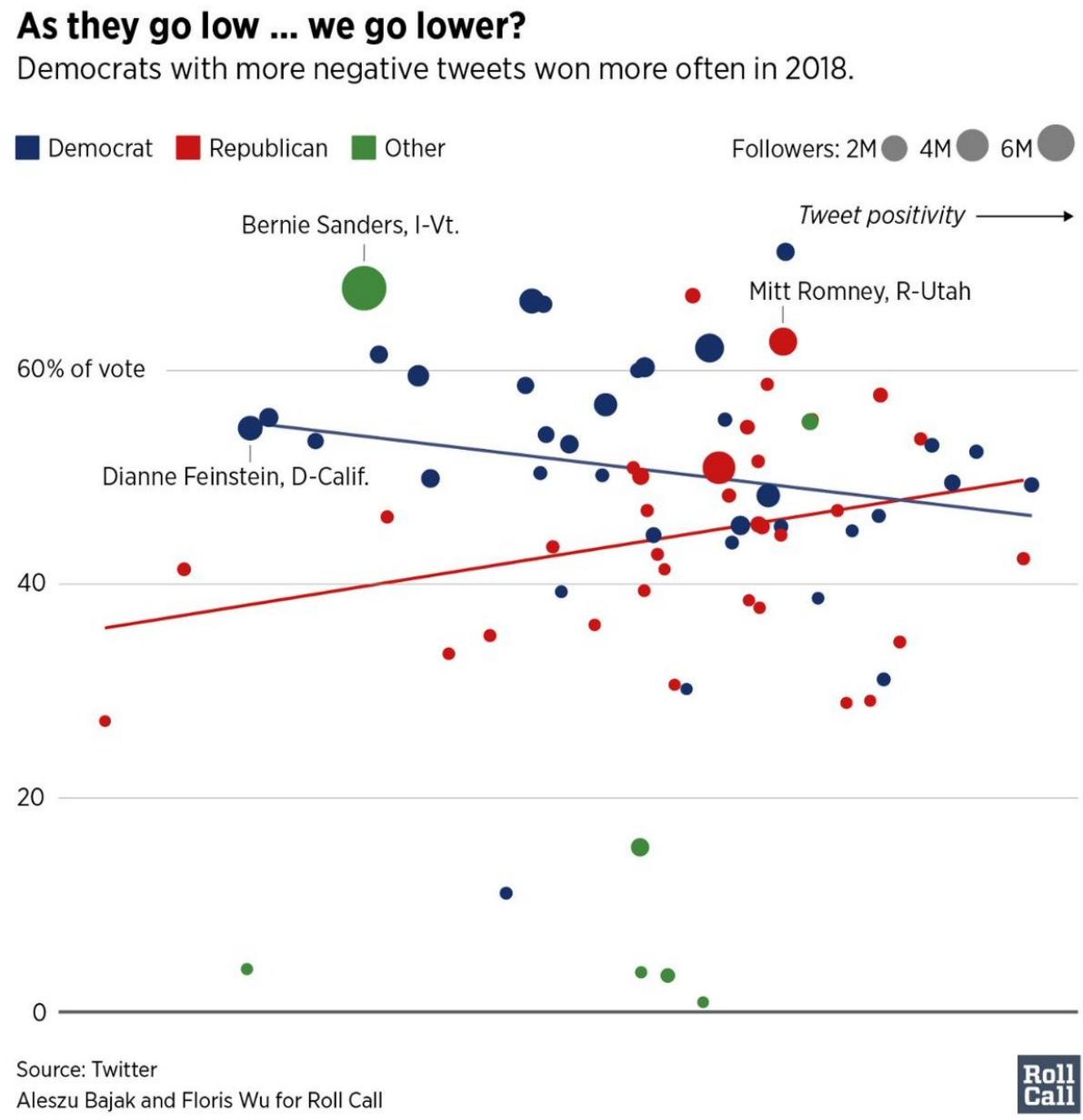
REPUBLICAN  
≈ right



# US Democrats vs US Republicans

- Studies based on volume and sentiment of postings
- Automated with tools from social media monitoring companies
- Claims often overstretched by journalists, authors, or both

*Tip: understate results for journalists, they will exaggerate them anyway*



# (Not) predicting election results

- A method of prediction should be an algorithm finalized **before** the election
  - specify data collection, cleaning, analysis, interpretation...
- Data from social media are fundamentally different than data from natural phenomena
  - people change their behavior next time around
  - spammers & activists will try to take advantage
- From a testable theory on *why* and *when* it predicts (avoid self-deception!)
- (maybe) Learn from professional pollsters
  - tweet ≠ user
  - user ≠ eligible voter
  - eligible voter ≠ voter

# Predict political stance

Results using an SVM

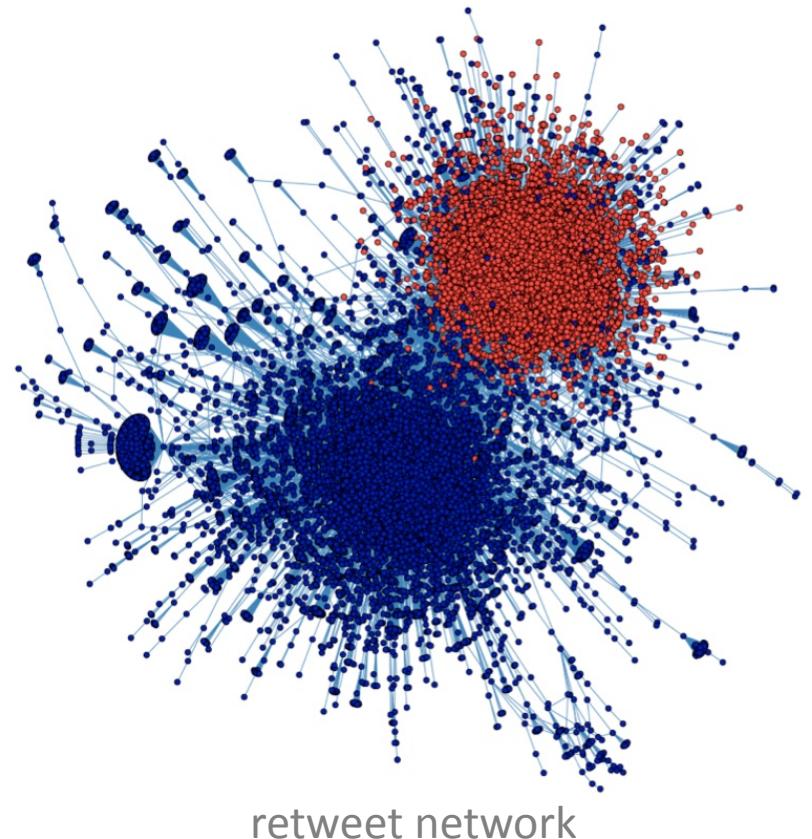
Features	Conf. matrix	Accuracy	
Full-Text	$\begin{bmatrix} 266 & 107 \\ 75 & 431 \end{bmatrix}$	79.2%	
Hashtags	$\begin{bmatrix} 331 & 42 \\ 41 & 465 \end{bmatrix}$	90.8%	
Clusters	$\begin{bmatrix} 367 & 6 \\ 38 & 468 \end{bmatrix}$	94.9%	network-based method
Clusters + Tags	$\begin{bmatrix} 366 & 7 \\ 38 & 468 \end{bmatrix}$	94.9%	

# Network-based methods

- **Label propagation**
  - Random initialization of labels
  - Update according to majority label among neighbors, break ties arbitrarily
- Use known node labels to determine which cluster means what

Network	Min	Max	Mean
Mention	0.80	1.0	0.89
Retweet	0.94	0.98	0.96

Adjusted Rand Index for 100 label propagation runs on political data



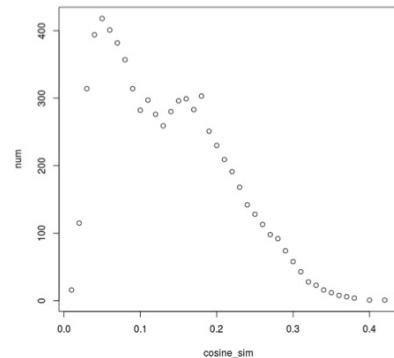
retweet network

Clusters  $\begin{bmatrix} 367 & 6 \\ 38 & 468 \end{bmatrix}$  94.9%

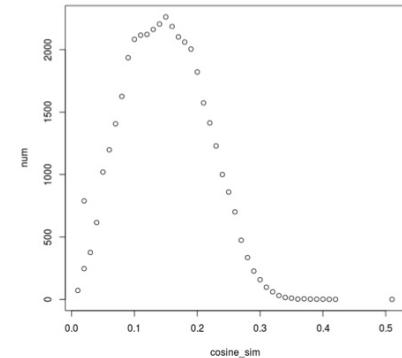
class assignment by cluster majority

# Secular vs Islamist

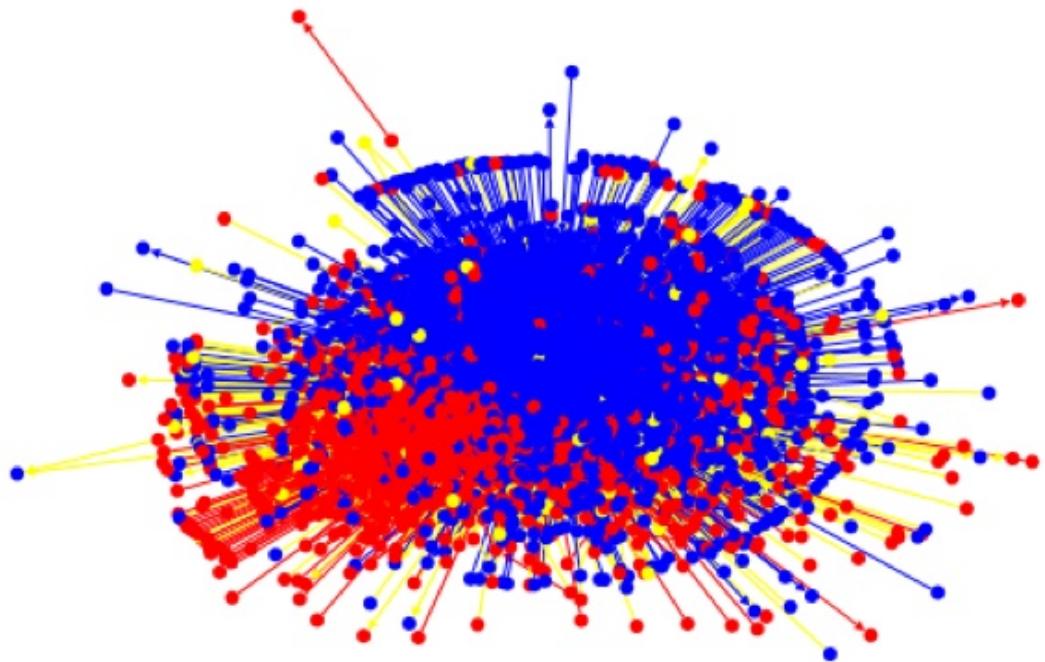
hashtag cosine user-user similarity

**Islamist**

average: 0.16  
median: 0.11

**Secularist**

average: 0.16  
median: 0.21



**Islamist, Secular, intra-ideology**

the closer to Islamist, use of  
 –religious terms *increases*  
 –charity-related terms *increases*  
 –derogatory terms *decreases*

# Exposure of Chinese tourists to protests in Hong Kong

- **Study group**

Chinese who traveled to Hong Kong before the protests started

- **Control**

returned to China 36 to 6 days before the protests started

- **Treatment**

returned after the protests, and hence possibly witnessed them

[Zhang 2015]

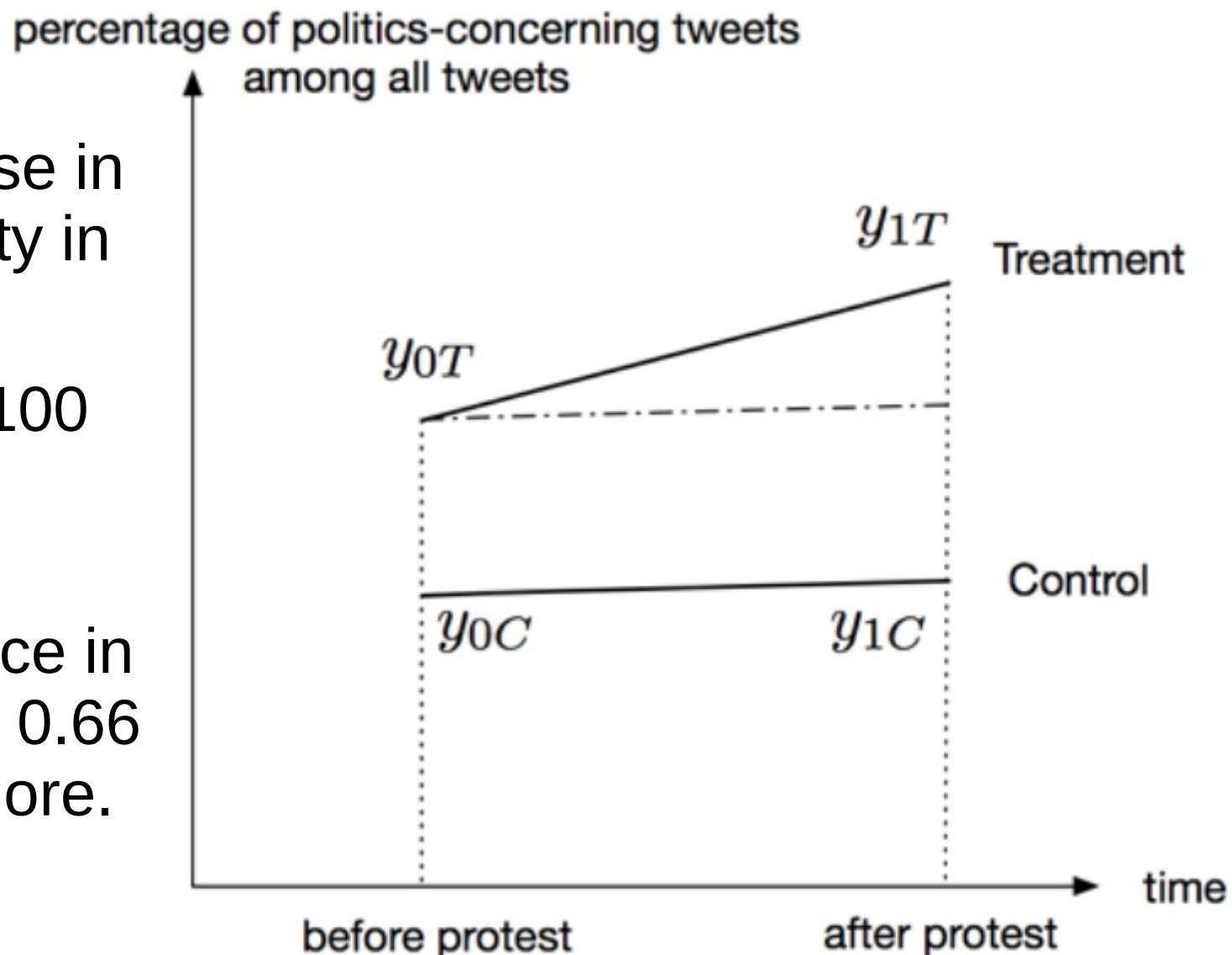


# Treatment vs matched control group

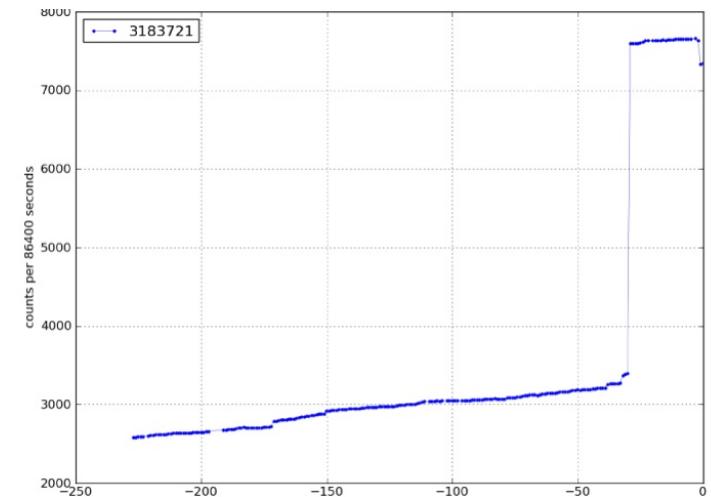
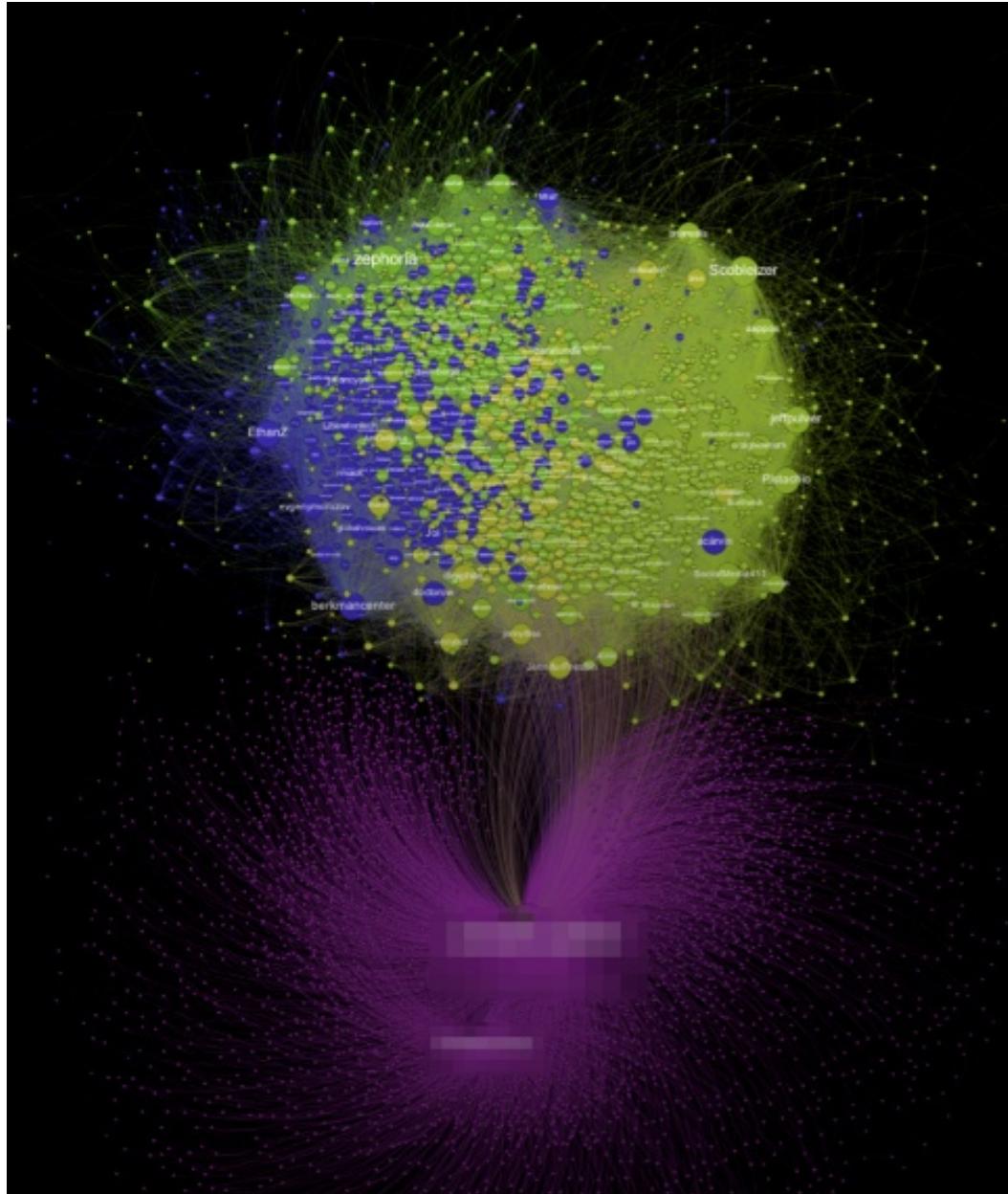
	<i>T</i>	<i>C</i>	<i>C<sub>match</sub></i>
Size	355	10169	710
Mean number of posts	787.30	726.96	772.13
Mean number of check-in	68.36	72.02	64.09
<b>Gender</b> =Male	108	3283	222
Female	176	6154	418
M/F Ratio	61.36%	53.35%	53.11%
Not Self-labelled	49	742	70
<b>Education</b> =College and beyond	108	3849	231
Not Self-labelled	225	6320	479
<b>Age</b> < 25	68	1806	129
<b>≥25</b>	8	163	14
Not Self-labelled	257	8200	567

# Difference-in-differences

- Found increase in political activity in Weibo
- Before: 1.66/100 posts were political posts
- After: difference in differences of 0.66 posts: 75% more.



# 4K followers x \$5



Faster organic followers  
acquisition after paid  
followers acquisition

**Fake Friends with  
Real Benefits  
by Gilad Lotan**

# Political Spam (“Truthy”)

Classifier	Resampling?	Accuracy	AUC
AdaBoost	No	92.6%	0.91
AdaBoost	Yes	96.4%	0.99
SVM	No	88.3%	0.77
SVM	Yes	95.6%	0.95

Ratkiewicz, J., Conover, M. D., Meiss, M., Gonçalves, B., Flammini, A., & Menczer, F. M. (2011, July). Detecting and tracking political abuse in social media. In Fifth international AAAI conference on weblogs and social media.

nodes	Number of nodes
edges	Number of edges
mean_k	Mean degree
mean_s	Mean strength
mean_w	Mean edge weight in largest connected component
max_k(i,o)	Maximum (in,out)-degree
max_k(i,o)_user	User with max. (in,out)-degree
max_s(i,o)	Maximum (in,out)-strength
max_s(i,o)_user	User with max. (in,out)-strength
std_k(i,o)	Std. dev. of (in,out)-degree
std_s(i,o)	Std. dev. of (in,out)-strength
skew_k(i,o)	Skew of (in,out)-degree distribution
skew_s(i,o)	Skew of (in,out)-strength distribution
mean_cc	Mean size of connected components
max_cc	Size of largest connected component
entry_nodes	Number of unique injections
numTruthy	Number of times ‘truthy’ button was clicked
sentiment scores	Six GPOMS sentiment dimensions



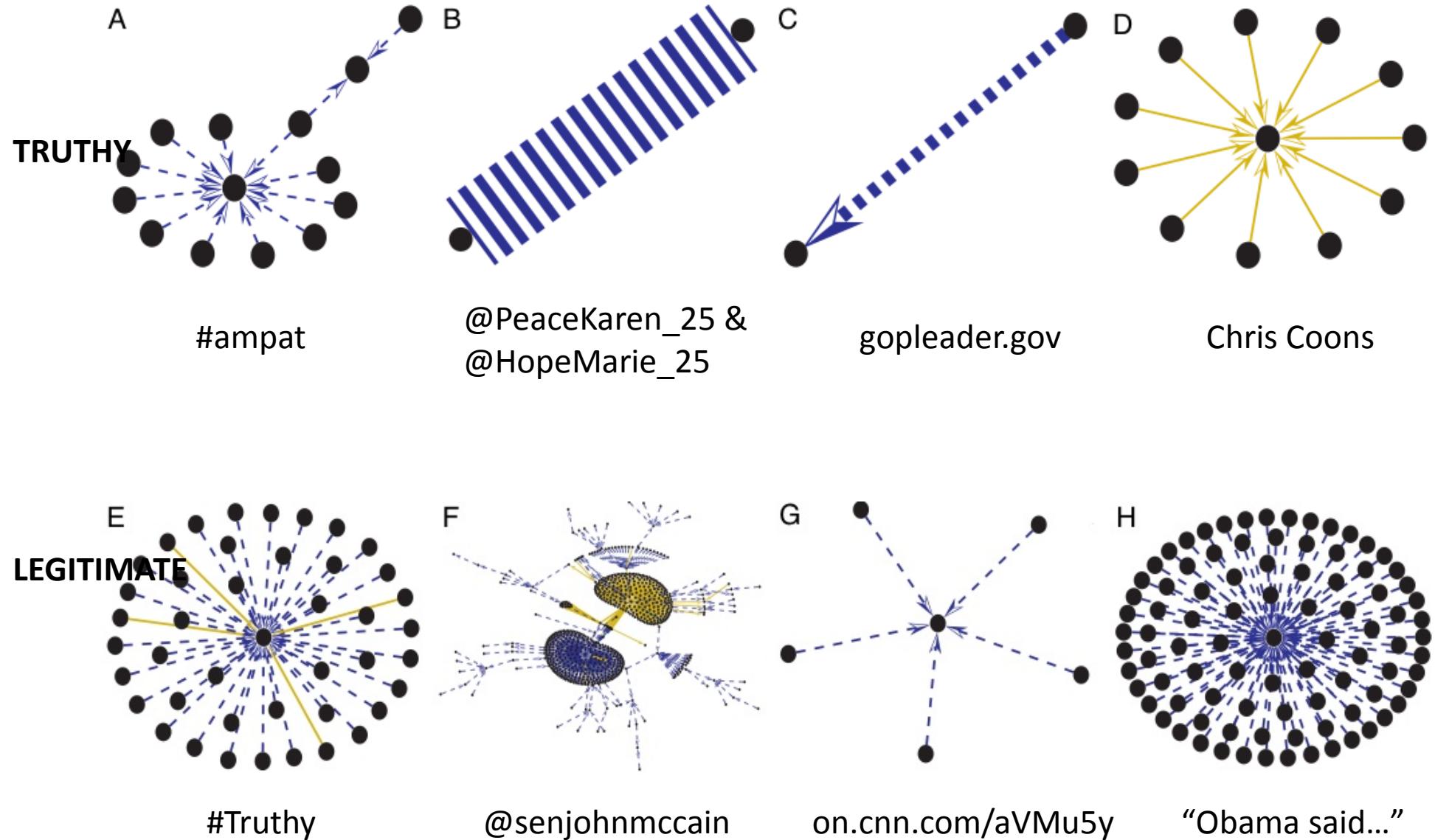
**Truthiness** is a quality characterizing a "truth" that a person making an argument or assertion claims to know intuitively "from the gut" or because it "feels right" without regard to evidence, logic, intellectual examination, or facts.

## Social bots distort the 2016 U.S. Presidential election online discussion

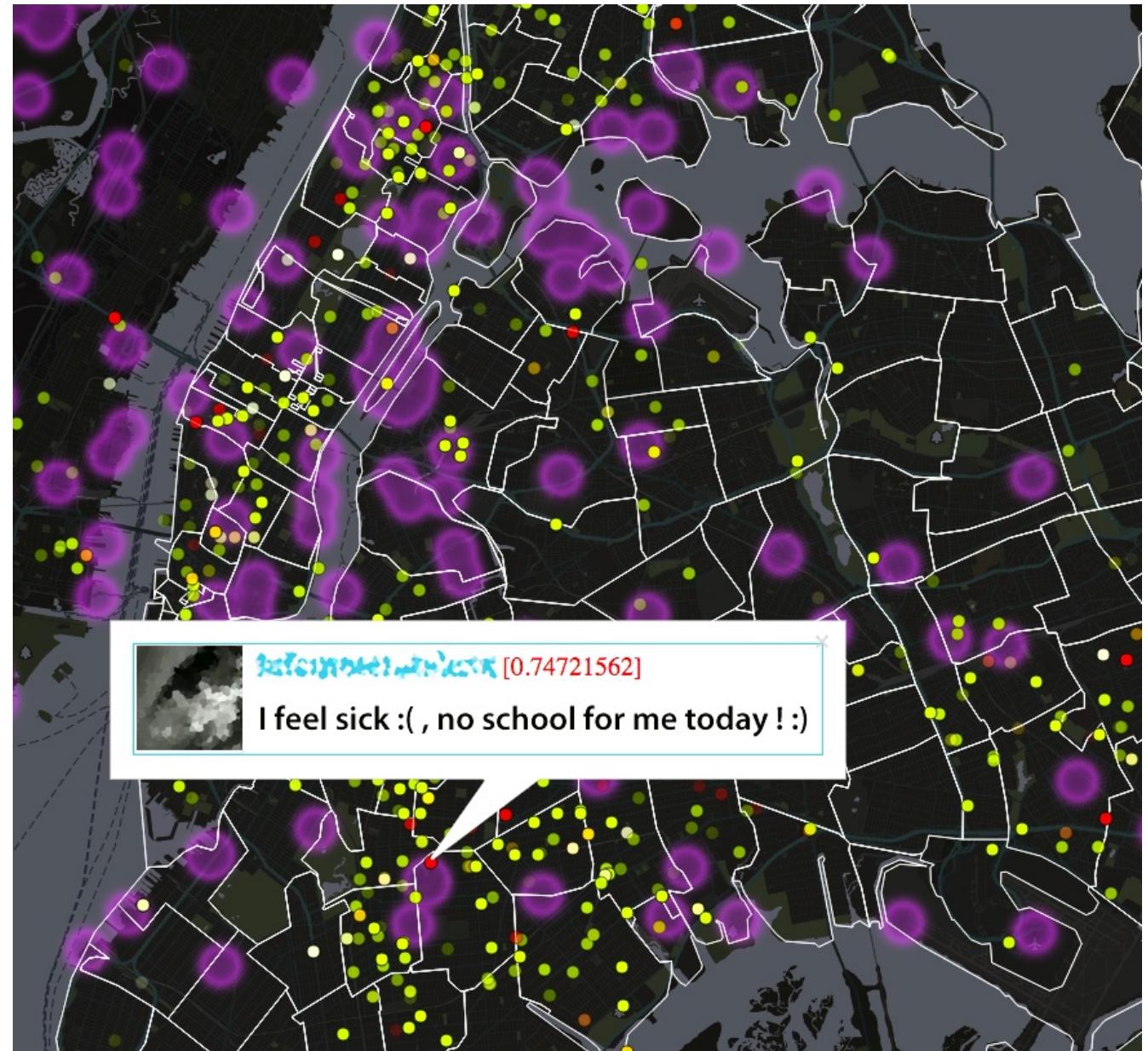
Bot specific statistics	Top 50K users (exact)
# Bot-generated tweets	2,330,252 (18.45%)
# Human-generated tweets	10,303,251 (81.55%)
# Bots	7,183 (14.4%)
# Humans	40,163 (80.3%)
# Unknown	2,654 (5.3%)

Bessi, Alessandro, and Emilio Ferrara. "Social bots distort the 2016 US Presidential election online discussion." First Monday 21.11-7 (2016).

# (Ab)normal linking patterns



# Lifestyle and health at scale



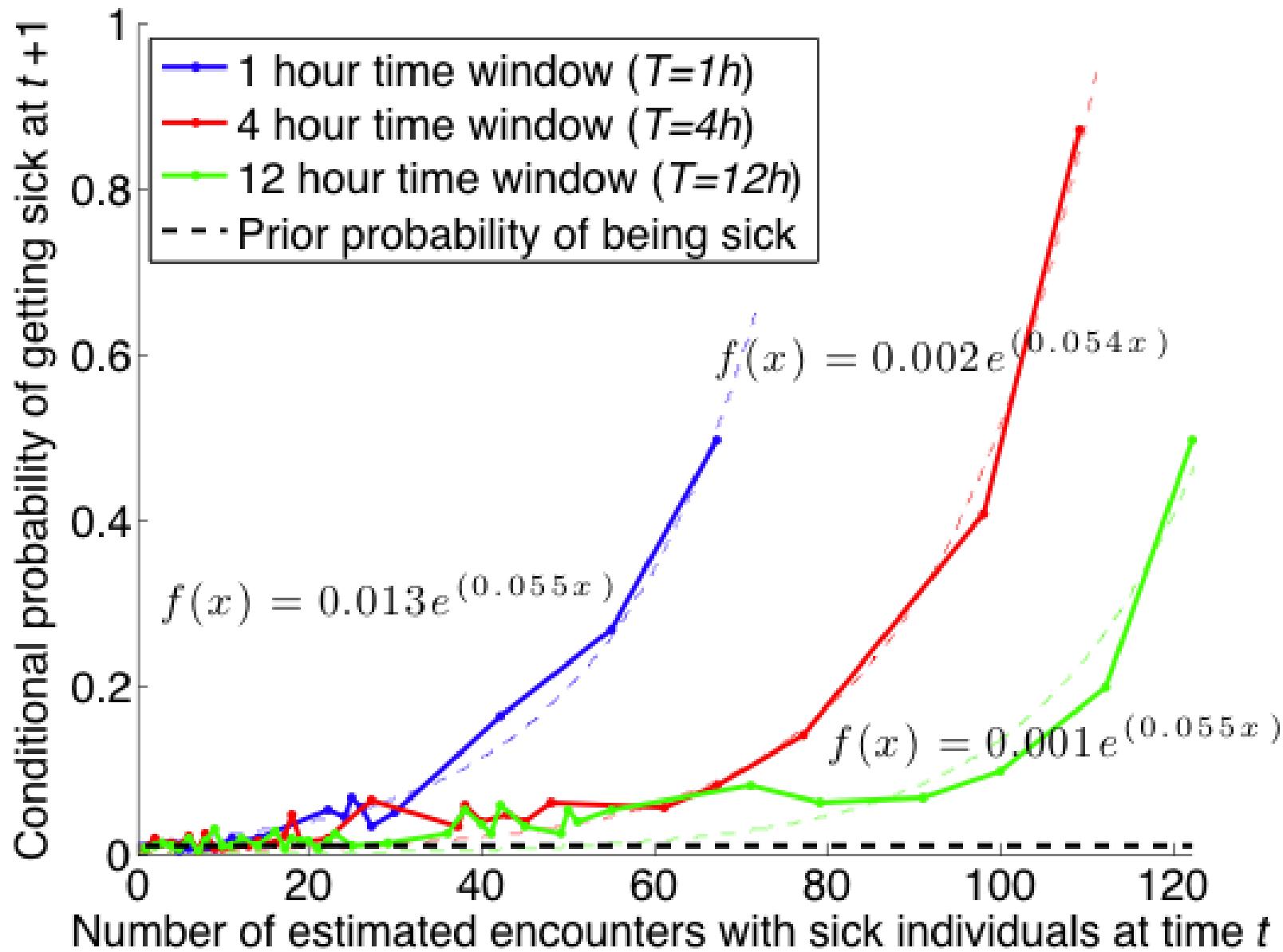
Sadilek et al.  
WSDM 2013

# Example features

Positive Features		Negative Features	
Feature	Weight	Feature	Weight
sick	0.9579	sick of	-0.4005
headache	0.5249	you	-0.3662
flu	0.5051	lol	-0.3017
fever	0.3879	love	-0.1753
feel	0.3451	i feel your	-0.1416
coughing	0.2917	so sick of	-0.0887
being sick	0.1919	bieber fever	-0.1026
better	0.1988	smoking	-0.0980
being	0.1943	i'm sick of	-0.0894
stomach	0.1703	pressure	-0.0837
and my	0.1687	massage	-0.0726
infection	0.1686	i love	-0.0719
morning	0.1647	pregnant	-0.0639

Positively and negatively weighted significant features.

# Probability of being sick increases with closeness to sick people



# Examples in “smart cities”

- Data-driven neighborhood boundaries
- Data-driven residential/commercial zones
- Tourism and beauty

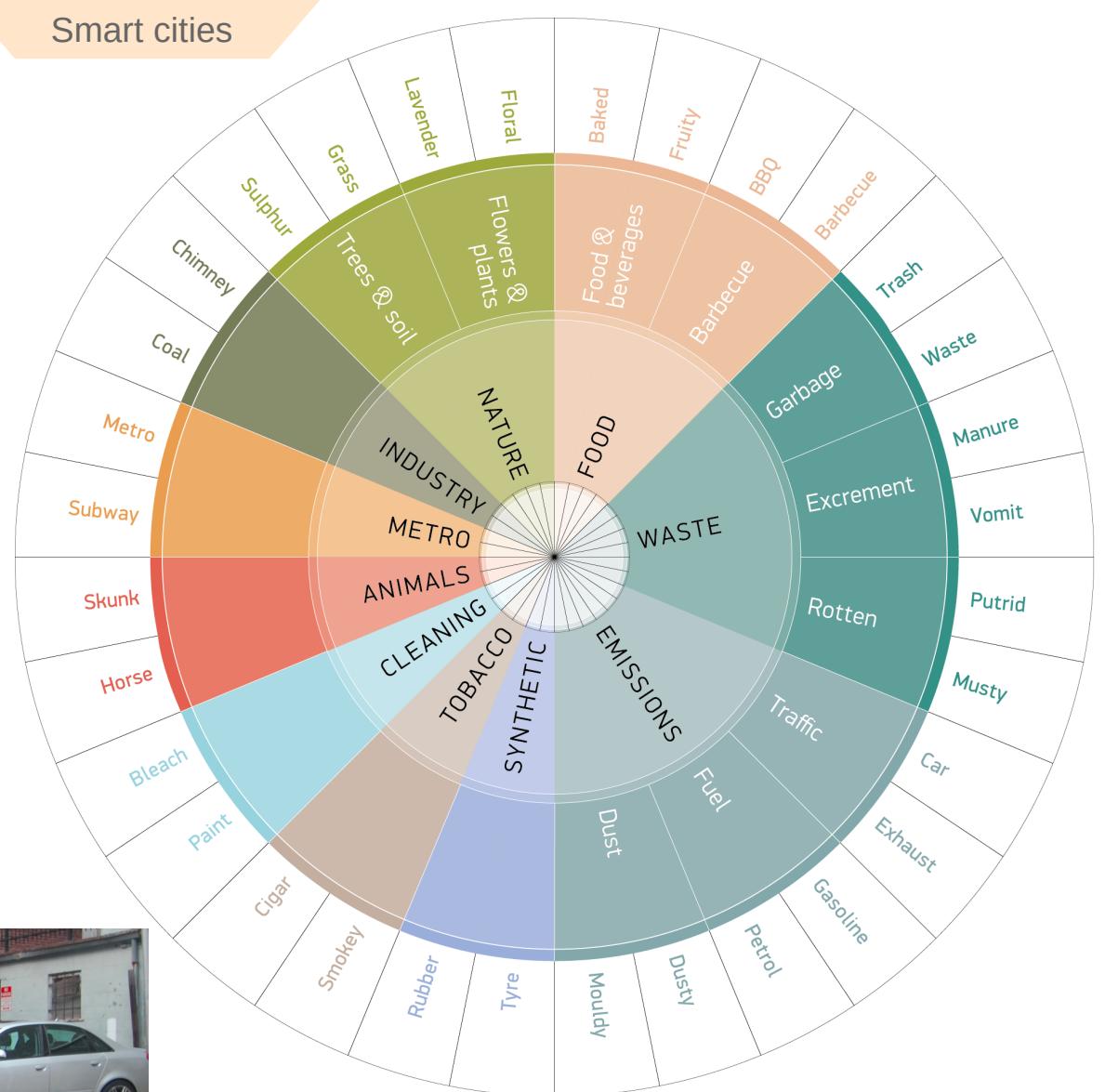
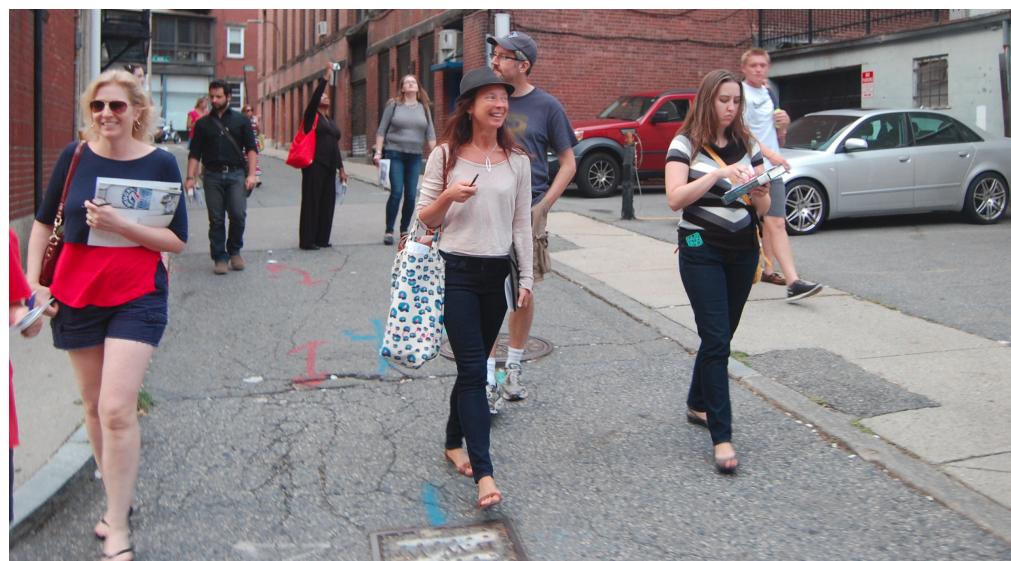


(a) Shortest



(b) Beauty

# Smell maps



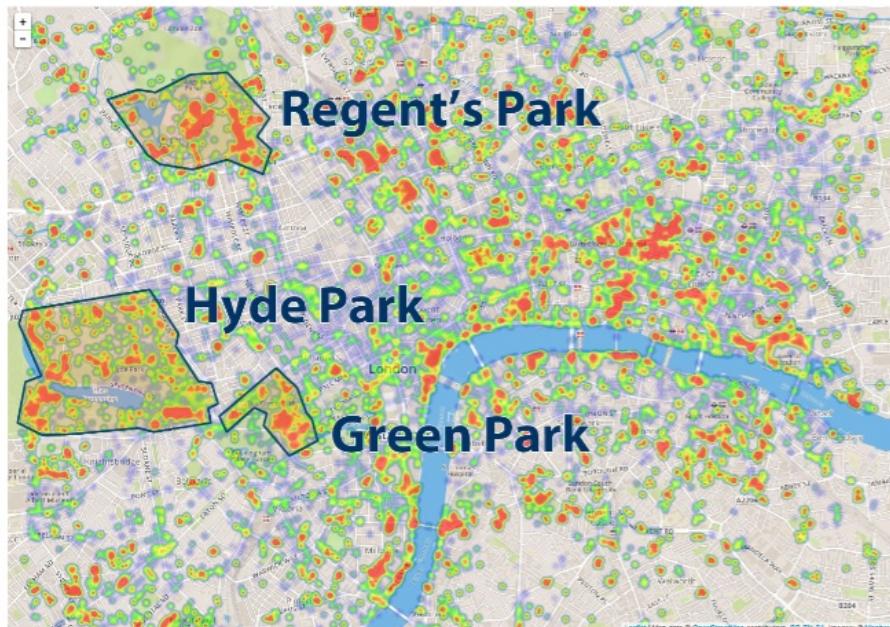
Urban Smellscape Aroma Wheel  
(depicting background and episodic aromas only)  
Aiello, L., McLean, K., Quercia, D., Schifanella, R., 2015

# Smells and tags

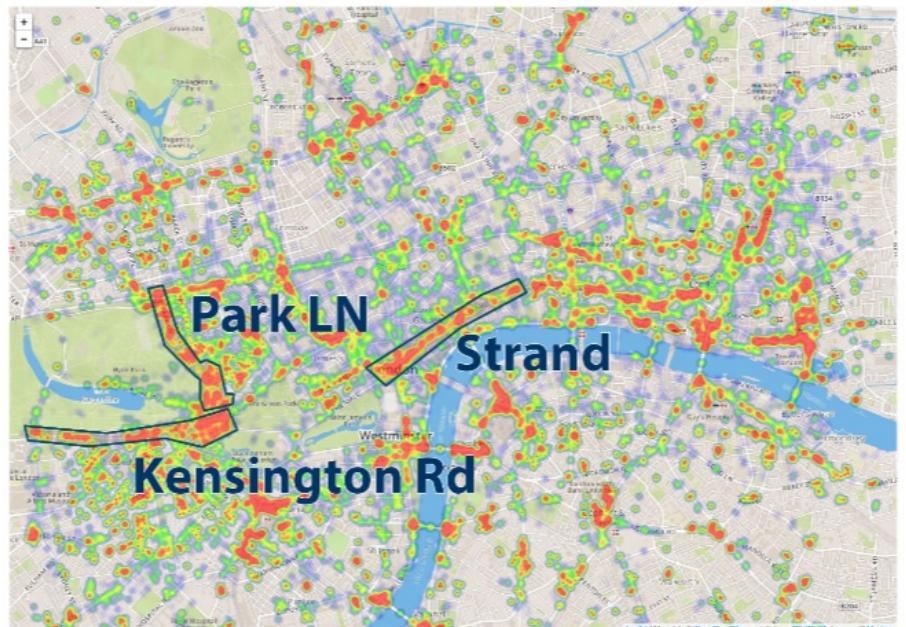
## Data: Instagram and Twitter

- Compute  $P(\text{smell}|\text{tags})$  for places where both
  - a set of tagged photos and
  - a set of smell annotations

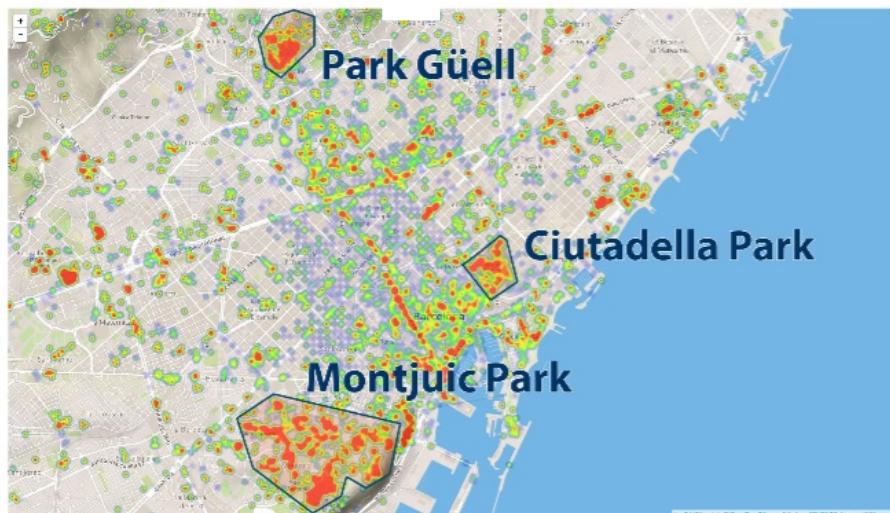
... have been observed
- Supervised learning approach



London, nature



London, emissions

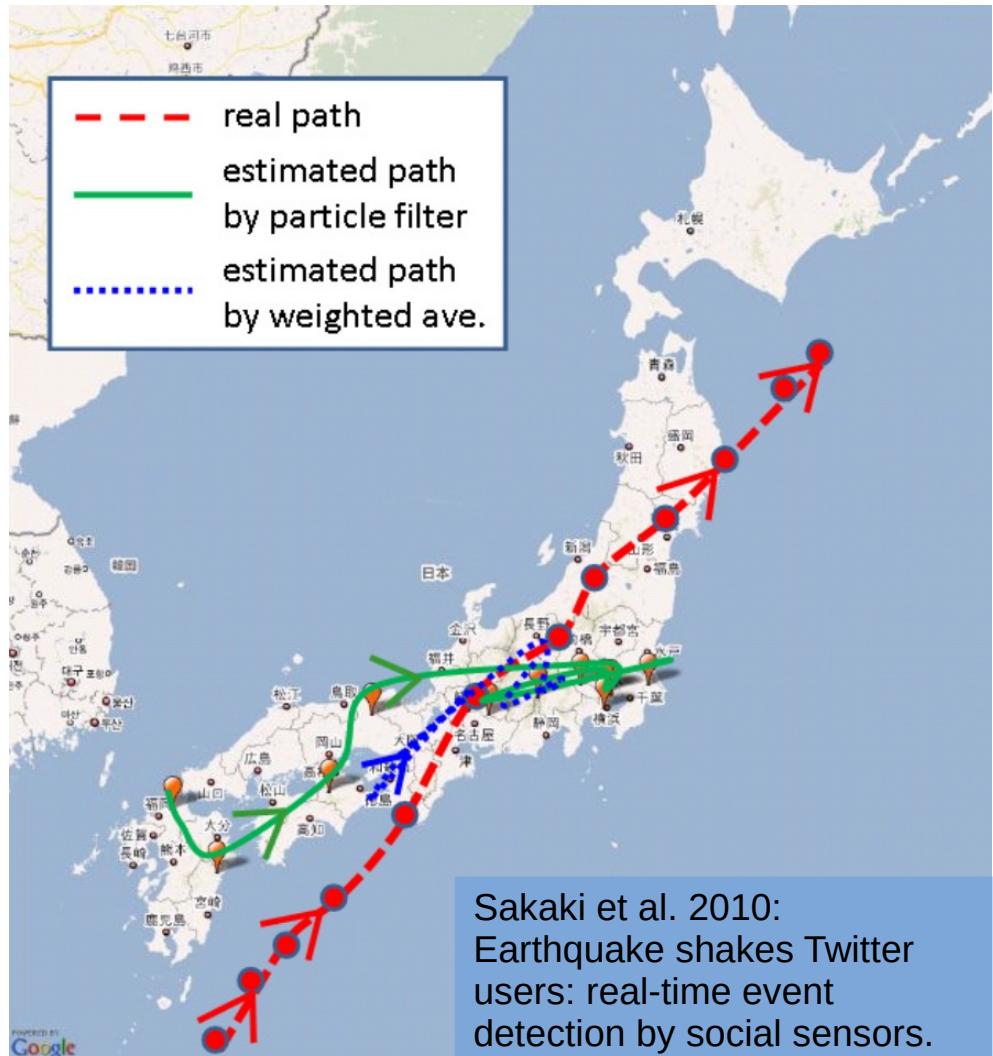


Barcelona, nature

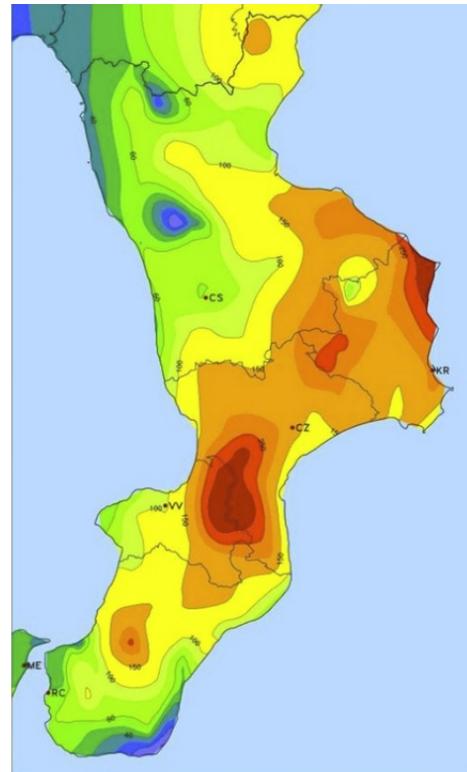


Barcelona, emissions

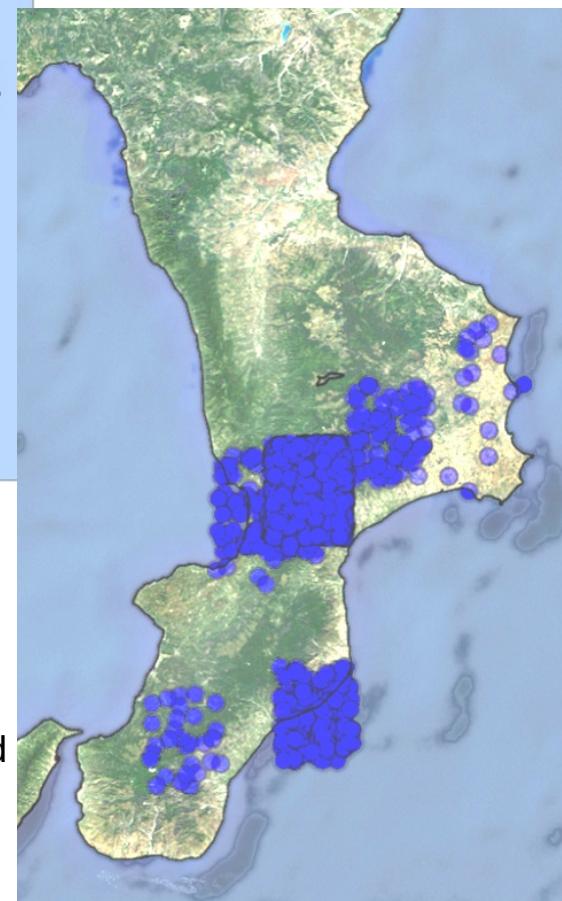
# Typhoon trajectory



# Flood risks

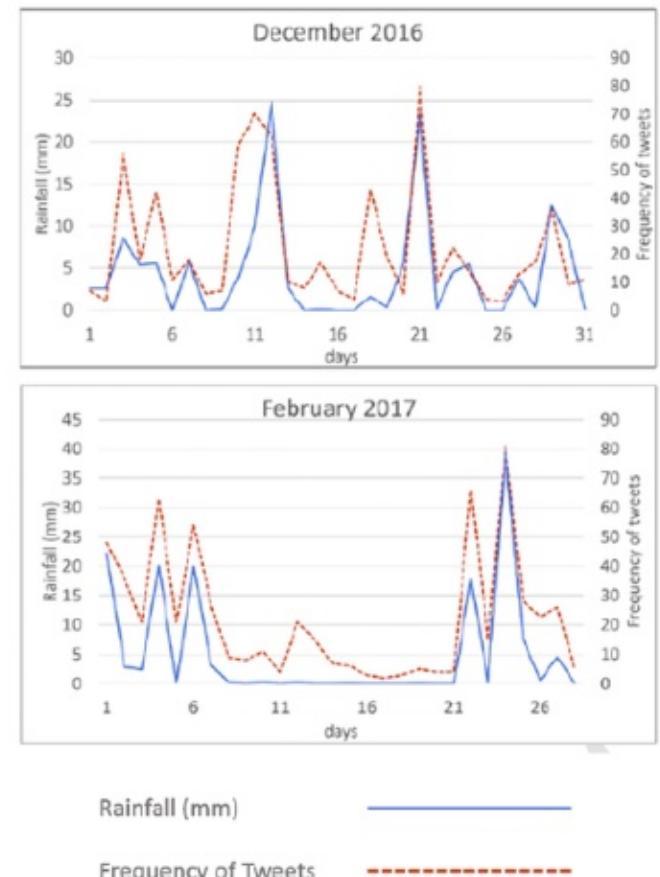
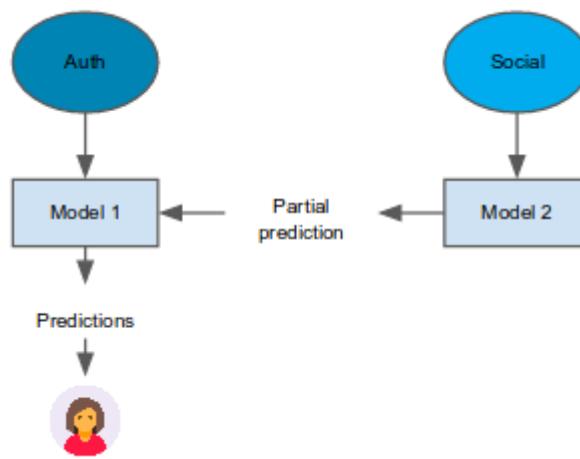


Lorini et al. 2019:  
Integrating Social Media into a Pan-European Flood Awareness System. Best paper award at ISCRAM.



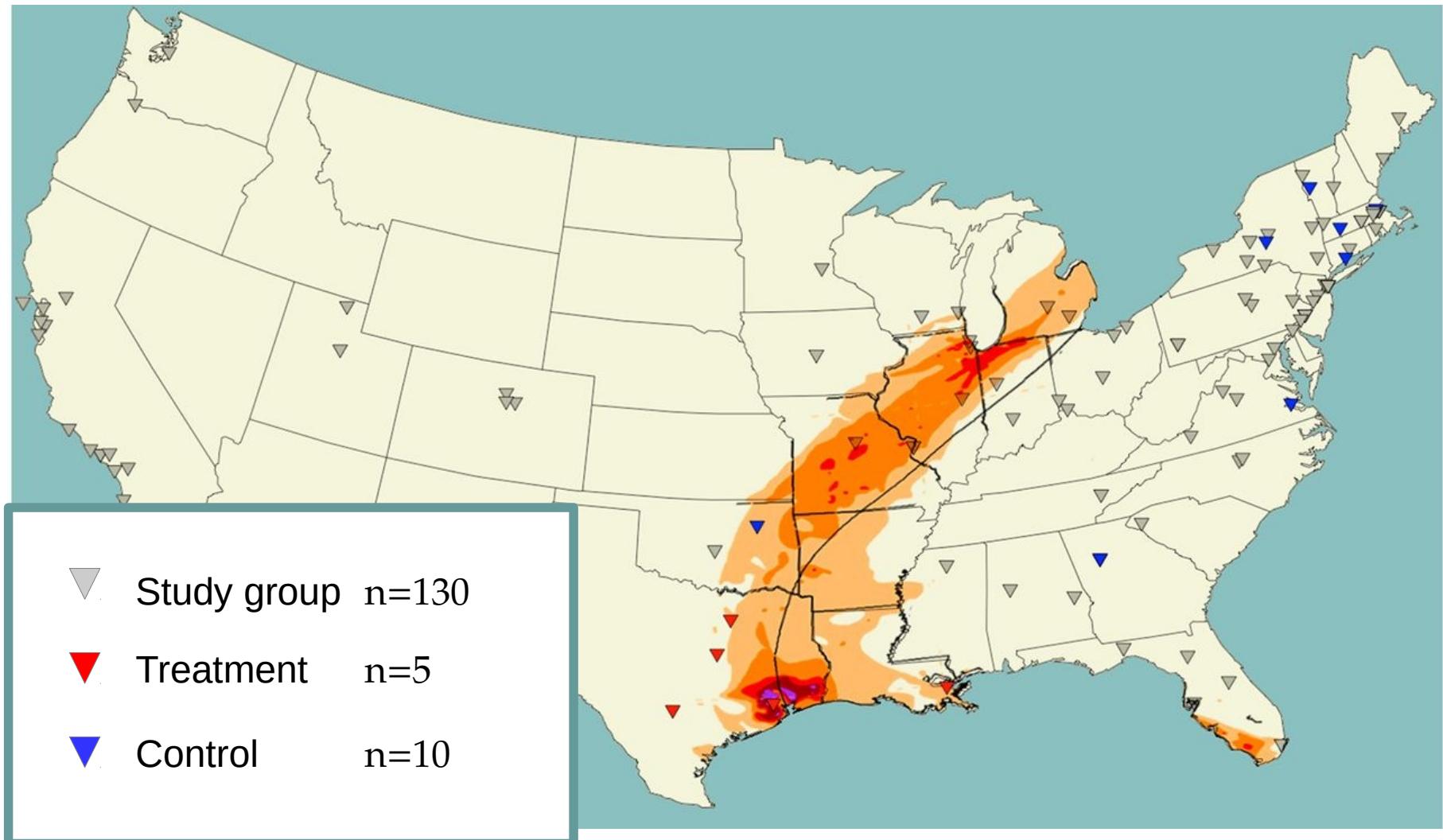
# From data points to the “Big Picture”

- Based on geocoded tweets (<2% of total) filtered to select those having keywords related to rain ("chuva" in Portuguese).
- Used to estimate rainfall (supervised setting), which is then fed into the flood risk prediction model



[Restrepo-Estrada et al. 2018]

# The effect of a hurricane in online social networks



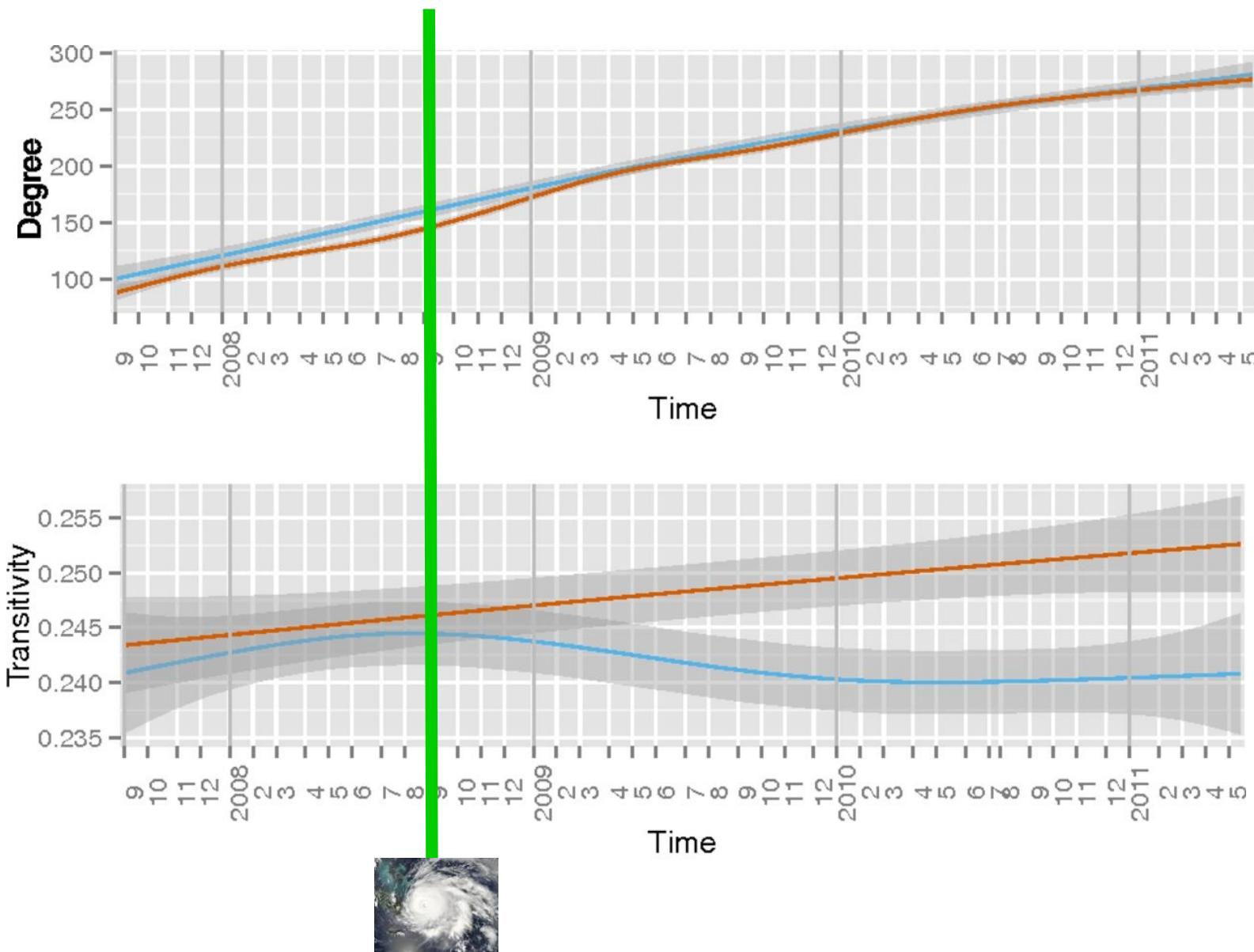
Phan, Tuan Q., and Edoardo M. Airoldi. "A natural experiment of social network formation and dynamics." *Proceedings of the National Academy of Sciences* 112.21 (2015): 6595-6600.

# Matching design: selection of control group

- Facebook posts from 1.5M students in 130 universities
- **Matched** 5 affected with 10 unaffected:
  - Similar: size, college ranking according to USNews, whether these colleges are public or private institutions, tuition fees, and other regional factors

Universities	No. of Users
Affected universities	
Baylor University	8,462
Rice University	2,355
Southern Methodist University	4,324
Trinity University	1,882
Tulane University	4,505
Unaffected universities	
Colgate University	2,359
The College of William and Mary	4,446
Georgia Institute of Technology	8,703
Middlebury College	2,374
Smith University	1,874
Tufts University	4,337
University of Pennsylvania	8,644
University of Tulsa	1,877
University of Utah	4,296
Yale University	4,519

# Results (red=treatment, blue=control)



Both undergo  
densification

Treatment has  
larger clustering  
coefficient  
(more triangles)

Beyond matching on specific attributes:  
**Propensity matching**

Alexandra Olteanu, Onur Varol, and Emre Kıcıman, Towards an Open-Domain Framework for Distilling the Outcomes of Personal Experiences from Social Media Timelines, in International Conference on Web and Social Media (ICWSM), AAAI - Association for the Advancement of Artificial Intelligence, 17 May 2016. [[link](#)]

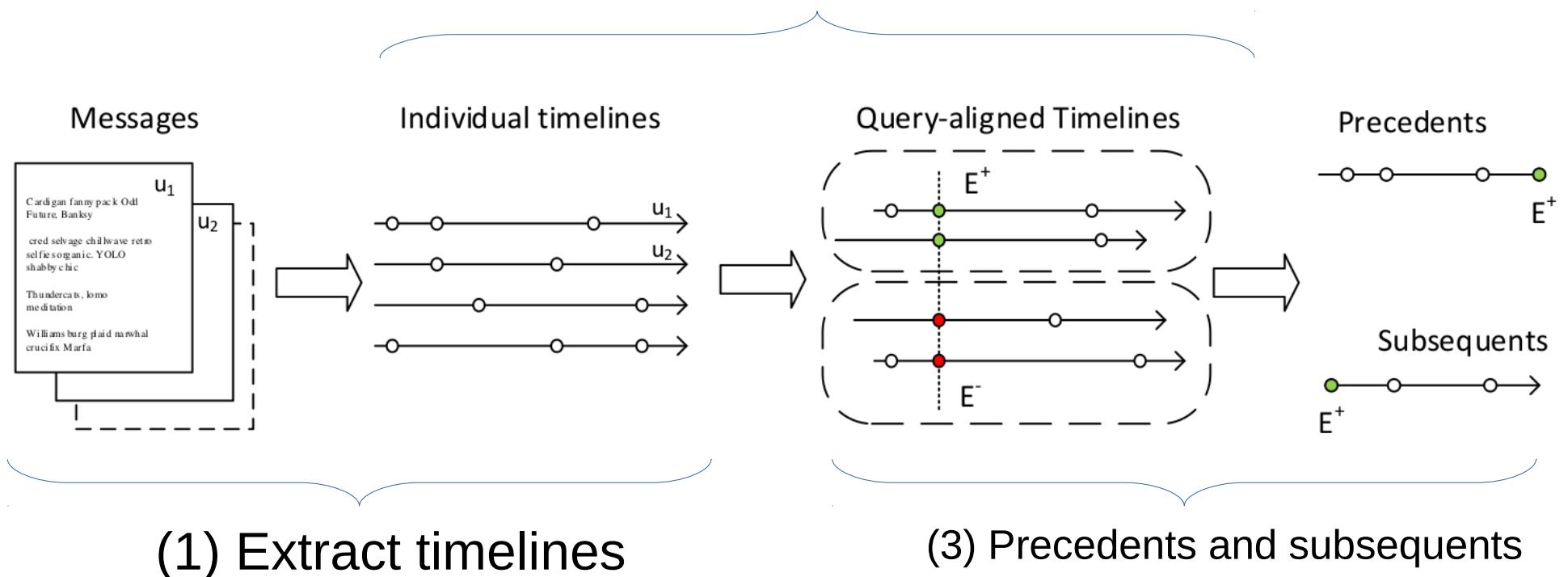
Idea: create a system to extract  
**Situation** → **Action** → **Outcomes**  
from social media

T<sub>1</sub>: “I got a kitten! We named  
her Versace :-)”

T<sub>2</sub>: “No sleep because the  
damn kitten is nuts!”

# Basic operations

## (2) Match events



# Many sub-problems

- Identification of experiential messages
- Timestamping event occurrences
- Recognition and canonicalization events
- Identification of precedent and subsequent events
- Identification of positive and negative valence of events

# Experiential messages classifier

Based on 10k labeled tweets.

**26% of tweets mention personal experiences**

8% mention goals/desires

66% are news/3<sup>rd</sup> person or other tweets

Personal Experiences	Other (news, 3 <sup>rd</sup> person, etc.)
Just completed a 15.72km run with @RunKeeper. Check it out! <URL> #RunKeeper	New campaign to protect children from second hand smoke launched <URL>
Just to set the mood I brought some Marvin Gaye and Chardonnay	Whoa. The kid from Cincinnati just suffered a horrible injury. Not good.
Lacrosse is so much fun why didn't I start earlier lol	@Bob I hear you.
Oh yeah guys we got a new puppy	@Charlie did you enjoy your night at the club?

# Event identification

I got a new kitten and he has blue eyes and stripes and I need a good name but nothing that's normal

Kıcıman, Emre, and Matthew Richardson. "Towards decision support and goal achievement: Identifying action-outcome relationships from social media." KDD 2015. [[link](#)]



I got a new kitten

== *got a cat,*  
*got a new cat,* ...



he has blue eyes



stripes



but nothing that's normal



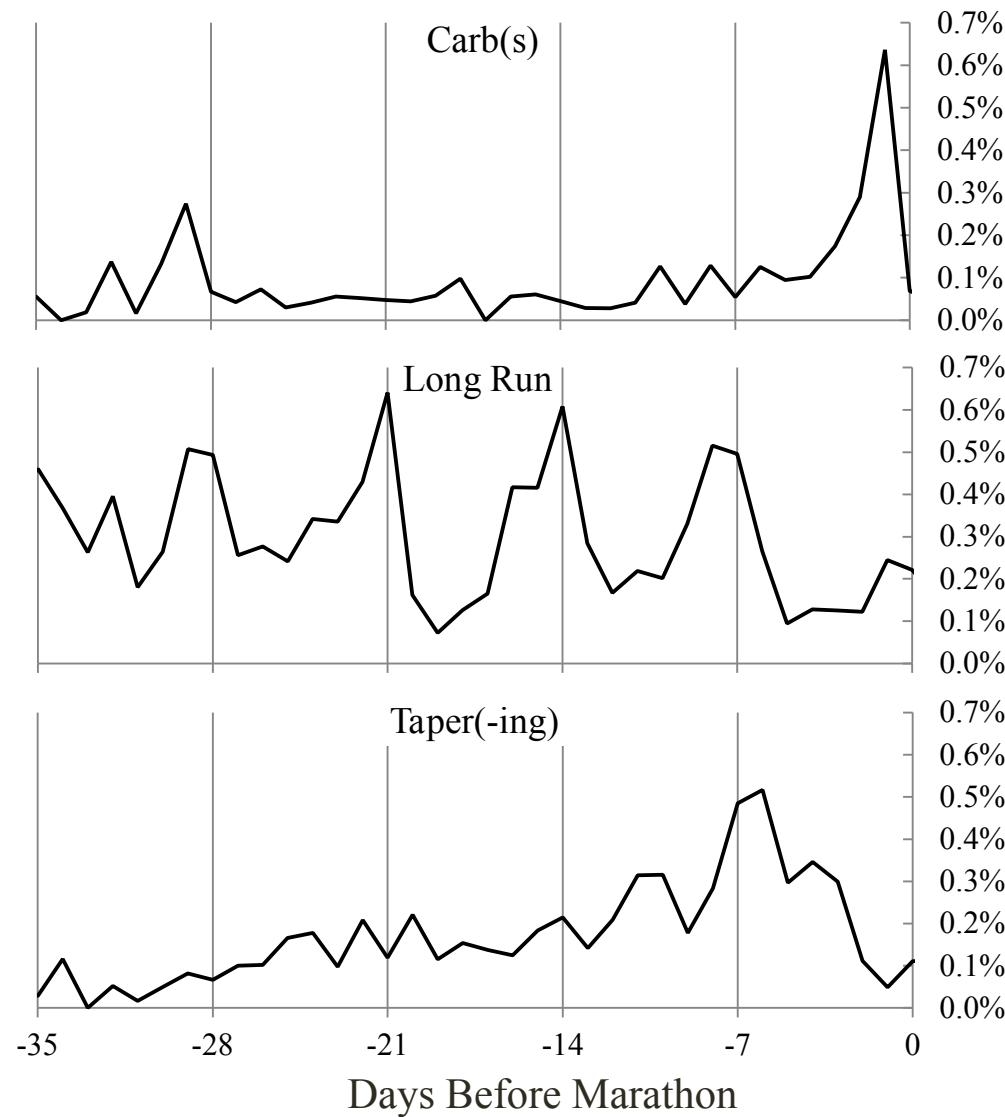
I need a good name

# Example subsequents

	Event	Example	PosNeg
Pros	cat named	We just got a cat and named it Versace	0.70
	I've got a cat	I've got a kitten asleep on my lap, and my heart has softened.	0.67
	Love my new kitten	I love my new kitten	0.88
Cons	Ran upstairs	But I ran upstairs and fell and now my head hurts	0.20
	Damn kitten	... no sleep because the damn kitten kept going nuts...	0.22
	Cat is literally	My cat is literally the devil	0.31

# Example precedents

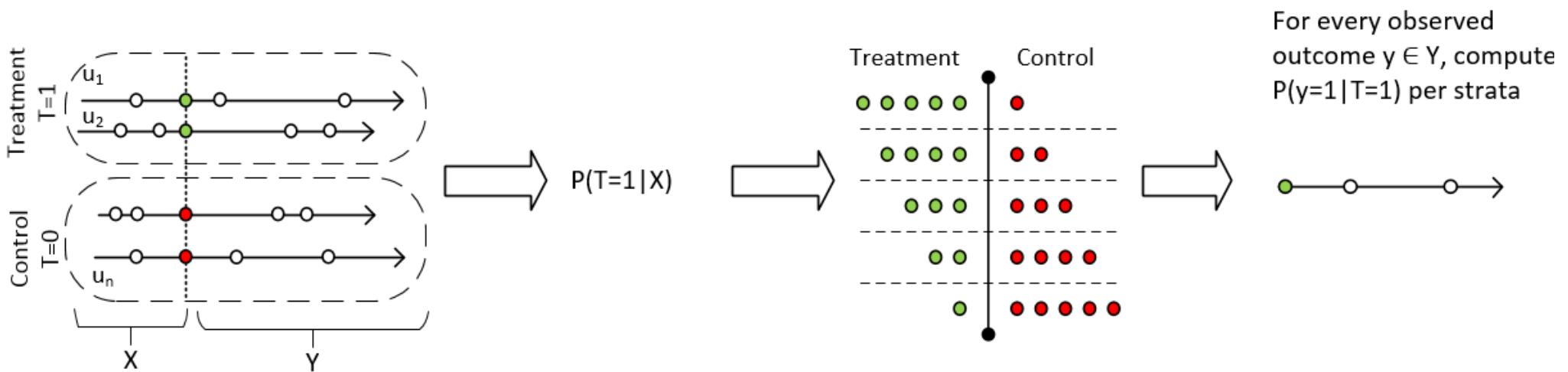
- Event: “personal record” in marathon



# Propensity matching explained

- You got a kitten
- According to what's known about your past, your probability of getting a kitten was  $x$
- You will be matched with someone whose probability of getting a kitten was also  $x$ 
  - But who did not get a kitten
- Every strata has a different unbalance
  - Which is predictable

# Propensity matching



Features of a user are all of their past events

PS Estimator trained w/average perceptron learning algorithm; extracted timelines are training data.

Decile stratification

# Matching design

- 39 situations in 9 groups
- Outcome is binary variable
- Average effect

$$P(\text{outcome}|T) - P(\text{outcome}|C)$$

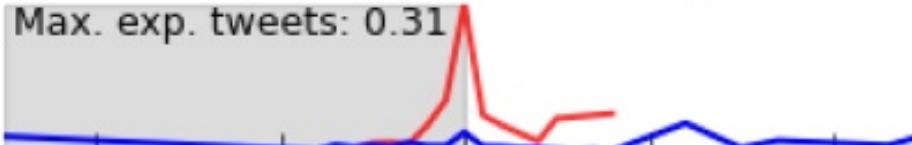
Category	Event	Treatment		
		Users	Msgs	
Business	Construct. and Maintenan.	Building stairs Cleaning countertops Installing a garbage disposal Painting the deck	24 8 29 592	8.3K 3.7K 8.2K 164K
	Financial Services	Owning a good credit card Paying credit card debts Buying life insurance Having pension	2291 414 1881 2344	920K 233K 561K 796K
		Incorporating one's business	28	9.1K
		Becoming a broker Investing money	855 23981	355K 9.8M
			Total	32447 12.8M
Health	Diseases	Having high blood pressure Having gout Having high cholesterol Having kidney stone Having high triglycerides lev.	5279 364 1384 727 27	1.9M 118K 522K 259K 12K
		Suffering from depression Suffering from OCD Being a sociopath Being a psychopath Suffering from anxiety	25207 11429 1491 2895 53983	10.5M 4.8M 676K 1.3M 22.6M
		Suffering from bipolar disord.	13723	6.3M
	Pharmacy	Taking Prozac Taking Lorazepam Taking Promethazine Taking Tramadol Taking Xanax	617 47 242 397 3300	222K 19.4K 118K 161K 1.4M
		Total	121112	50.9M
Society	Issues	Losing belly fat Increasing gross income	93 135	24.8K 93.2K
		Getting divorced Becoming a notary	2717 65	1.2M 22.7K
	Law	Applying for social security Filing for bankruptcy Having a living trust	6172 921 18	2.3M 347K 7.5K
		Finding true love Recovering after adultery	1885 9	654K 4.4K
		Filing divorce Dealing with jealousy issues Changing last name	422 789 1019	178K 370K 403K
			Total	14245 5.7M

# Example: having high triglycerides level

Outcome	Count	Absolute Increase	Z-Score
Your_risk	46	24.8%	18.12
Statin	48	23.1%	17.69
Lower	120	35.9%	17.18
Cardiovascular	54	23.0%	16.72
Healthy_diet	55	19.3%	16.54
Fatty_acid	29	18.3%	16.37
Help_prevent	73	26.9%	16.01
Risk_factor	33	18.3%	15.55
Fish_oil	48	24.4%	15.42
inflammation	78	25.1%	15.30

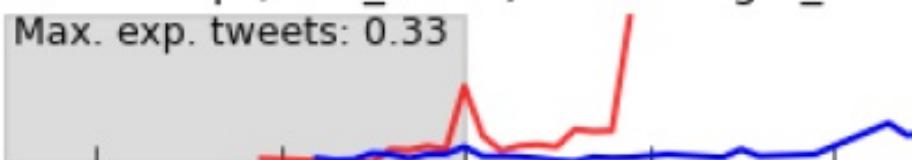
Pharmacy/Prozac, outcome:depression

Max. exp. tweets: 0.31



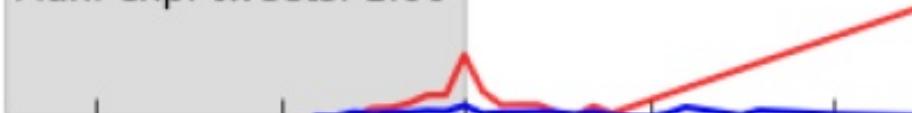
Relationships/Last\_name, outcome:get\_marry

Max. exp. tweets: 0.33



Relationships/Jealousy, outcome:jealous

Max. exp. tweets: 1.00



Pharmacy/Lorazepam, outcome:medication

Max. exp. tweets: 1.00



Pharmacy/Tramadol, outcome:painkiller

Max. exp. tweets: 0.12



Society/Belly\_fat, outcome:fitness

Max. exp. tweets: 0.91



Society/Gross\_income, outcome:get\_pay

Max. exp. tweets: 0.68



Diseases/Gout, outcome:joint

Max. exp. tweets: 0.24



Pharmacy/Tramadol, outcome:pain

Max. exp. tweets: 1.61



Pharmacy/Xanax, outcome:weed

Max. exp. tweets: 0.51



Treatment

Control

# Example Treatments and Outcomes (paraphrased for anonymity)

Treatment	Example tweet	Outcome	Example tweet
Dealing with jealousy issues	<i>ironically, u ask why I have jealousy issues</i>	wake up	<i>@user I need u to <u>wake up</u> because im bored</i>
Suffering from depression	<i>if u think depression is eccentric or cute u can have mine bc i dont wanna deal with it</i>	thoughts	<i>hate small talks, dont talk abt weather, tell me what keeps u up at night, ur thoughts abt dying</i>
Suffering from depression	<i>if u think depression is eccentric or cute u can have mine bc i dont wanna deal with it</i>	self harm	<i>@user dont self harm, remember yr worth so much better, u dont deserve this pain, stay safe</i>
Suffering from anxiety	<i>if hadnt spent years dealing with anxiety, I wouldnt have my sense of humor</i>	yelling	<i>dont have any anger issues at all, im really happy when <u>yelling</u> at people</i>
Paying credit card debts	<i>seriously, my soul was deep hurt when I paid that credit card bill</i>	apartment	<i>Im checking some <u>apartments</u> in NYC lol</i>
Having high blood pressure	<i>@user but I do have really <u>high blood pressure!</u></i>	stress	<i>had a heart mri and the news are as good as they can be. much <u>stress</u> is now removed from my life</i>
Having gout	<i>@user I know <u>I have gout...</u></i>	uric	<i>is the increase in <u>uric</u> acid production in blood, forming crystals and depositing them in joints</i>

Labeling by non-experts (Mechanical Turk workers)

Further reading: [observational studies and propensity matching](#)

Biases,  
biases  
all  
the  
way

Credit: [Sam Hollingsworth](#)



**Type I** research goals: understand/influence phenomena specific to social platforms

**Type II** research goals: understand/influence phenomena beyond social platforms

Construct validity

Internal validity

External validity



### General biases and issues

Population  
biases

Behavioral  
biases

Content  
biases

Linking  
biases

Temporal  
biases

Redundancy



### Biases at source

Functional biases  
Normative biases  
External biases  
Non-individuals

### Collecting

Acquiring  
Querying  
Filtering

### Processing

Cleaning  
Enriching  
Aggregating

### Analyzing

Qualitative analysis  
Descriptive statistics  
Inferences & predictions  
Observational studies

### Evaluating

Metrics  
Interpretations  
Disclaimers



Data platforms (not under researcher control)

Research designs (under researcher control)

# Take-home messages

# Best practices: social media mining

- **Interdisciplinary work**
  - Mixed methods: qualitative and quantitative
  - Robust to different settings, metrics, datasets
  - Aware that there are biases, biases all the way
- **Good applied research**
  - Is well-grounded in the literature of the domain
  - Has a measurable impact in the problem domain

# Practical advice: social media research



## This is a **very crowded area**

- Stand out with a completely new problem, a new dataset, and a deeper grounding on a problem domain
- Be the data you want to see in the world



## Become **learned** in your problem domain

- Commit to cross-disciplinary research or don't do it
- Talk to many experts on the problems you're trying to address; learn their language; read key books
- They'll dislike your results and/or jump to high-level conclusions that don't follow your low-level observations, probably both; listen and be patient