# Big Data Analytics

## FOSCA GIANNOTTI AND LUCA PAPPALARDO

# Suggested Bibliography

GoogleFluTrend

Measuring the Happiness of Large-Scale Written Expression: Songs, Blogs, and Presidents, Peter Sheridan Dodds Æ Christopher M. Danforth, J Happiness Stud (2010)

The Geography of Happiness: Connecting Twitter Sentiment and Expression, Demographics, and Objective Characteristics of Place Lewis Mitchell1 , Morgan R. Frank, Kameron Decker Harris1,2, Peter Sheridan Dodds, Christopher M. Danforth1
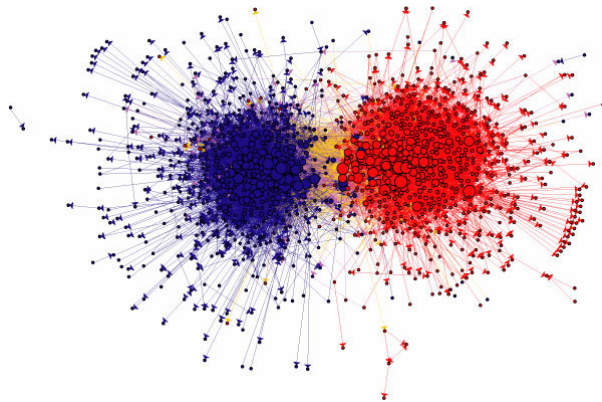
# SUMMARY

- Social media data and analytics (examples)

- Nowcasting Flu trend with search data –

- Predicting Happyness with lyrics ..and twitter – with predefined lexicon and mood weight

- Superdiversity

- Quantifying opinions with ISA algorithm – Forecasting Peruvian Elections

- Twitter data (UGC) for Migration Studies

# Social Media Data

Raw data +

Shares, Likes, Mentions, Impressions, Hashtag usage, URL clicks, Keyword analysis, New followers, Comments





Web 2 Icons

# WEB GALLERY OF ART & COLORS
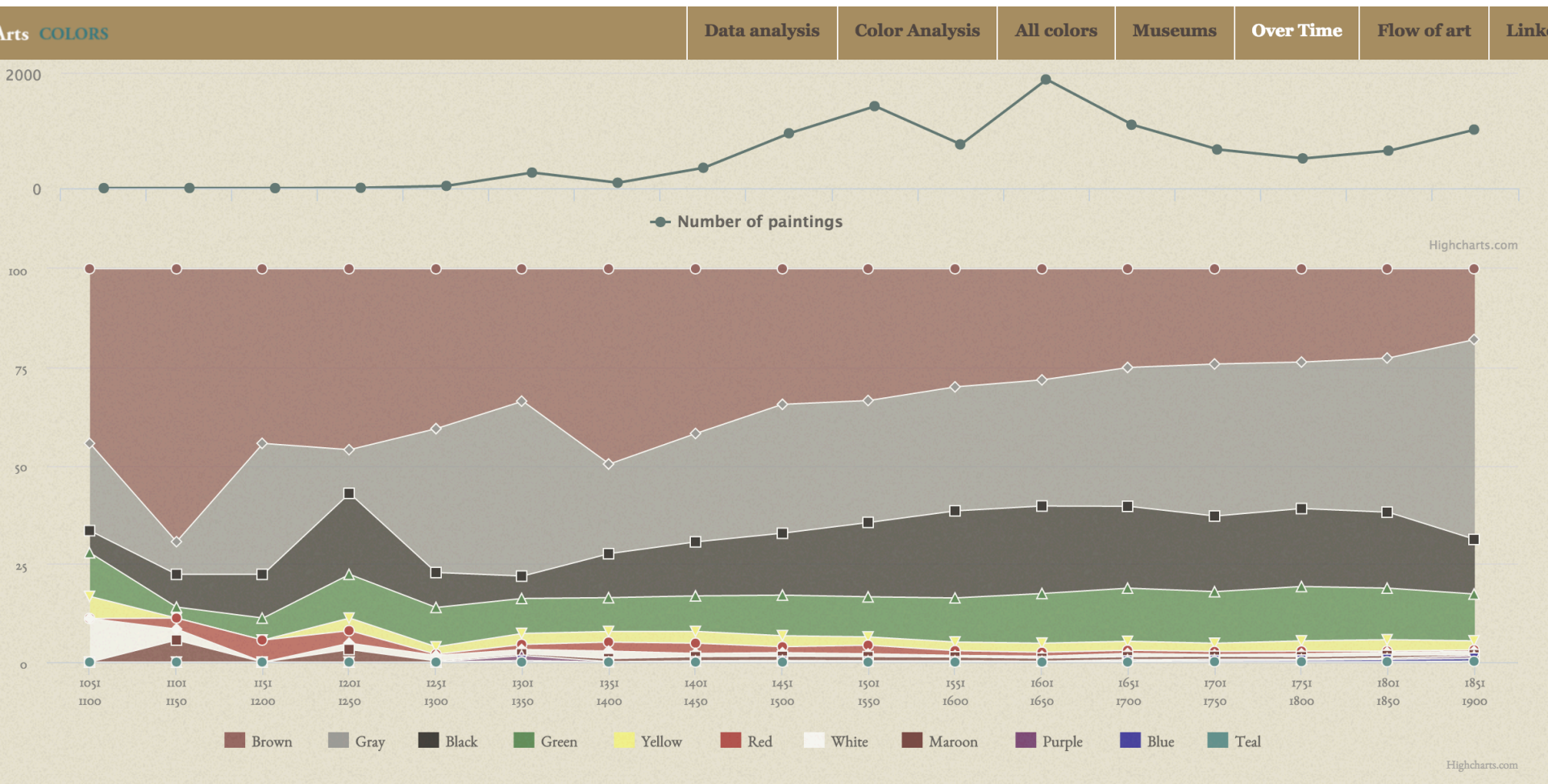
A visual representation of 850 years of paintings

CONCEPT    VISUALIZATION

http://sobigdata.danielefadda.com
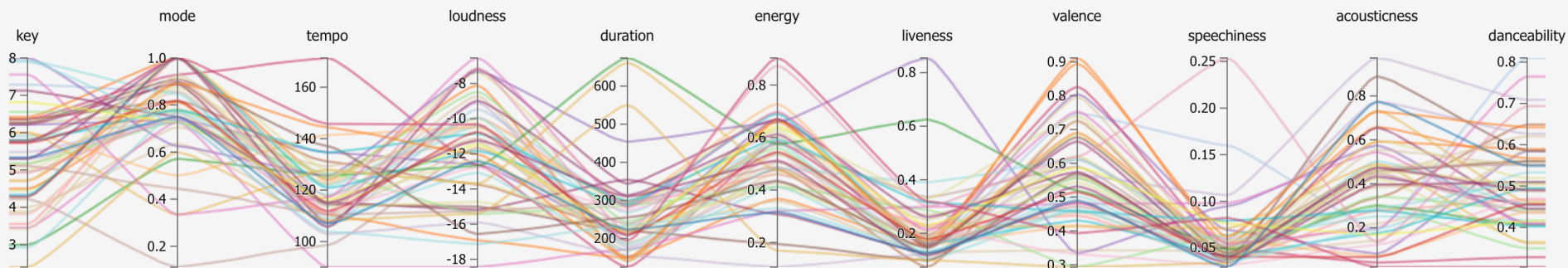
on exhibit at Data Stories 2015

# Rolling Stone 500 Albums — Top 50 Analysis

Keep | Exclude | Export

**Axes:** key | mode | tempo | loudness | duration | energy | liveness | valence | speechiness | acousticness | danceability

## The secret of success?

What you see is a multidimensional explorer of musical features from the top 50 Albums of Rolling Stone 500 Greatest Albums.

We used Echo Nest API and Python to get the musical features for every song in an album, then we cleaned, filtered and merged the files to get one JSON per album with the average features.

You can clearly see a common pattern between the albums; it seems that the "best" albums mostly share the same musical features.

## Controls

**Brush**: Drag vertically along an axis.
**Remove Brush**: Tap the axis background.
**Reorder Axes**: Drag a label horizontally.
**Invert Axis**: Tap an axis label.
**Remove Axis**: Drag axis label to the left edge.

## Credits

Visualization adapted from the Nutrient Database Explorer by Kai Chang @ 2012, all credits to him

Project Done by Samuele Borgheresi, Tommaso Ferrari Aggradi and Marco Da Campo

## Genres

- 1 Alternative Rock
- 2 Blues
- 1 Blues Rock
- 1 Country Rock
- 2 Folk
- 4 Folk Rock
- 1 Glam Rock
- 1 Grunge
- 2 Hard Rock
- 2 Jazz
- 1 Pop
- 6 Pop Rock
- 2 Progressive Rock
- 4 Psychedelic Rock
- 2 Punk
- 2 Punk Rock
- 1 Rap
- 1 Reggae
- 2 Rhythm and blues
- 5 Rock
- 3 Rock & Roll
- 1 Soft Rock
- 2 Soul
- 1 Southern Rock

## Sample of 25 Albums

Search…

- Bob Dylan - Blonde on Blonde (1966)
- Bob Dylan - Blood on the Tracks (1975)
- Bob Dylan - Bringing It All Back Home (1965)
- Bob Dylan - Highway 61 Revisited (1965)
- Carole King - Tapestry (1971)
- Chuck Berry - The Great Twenty-Eight (1982)
- David Bowie - The Rise and Fall of Ziggy Stardust and the Spiders from
- James Brown - Live at the Apollo (1963)
- John Coltrane - A Love Supreme (1965)
- Led Zeppelin - Led Zeppelin (1969)
- Love - Forever Changes (1967)
- Michael Jackson - Thriller (1982)
- Miles Davis - Kind of Blue (1959)
- Pink Floyd - Dark Side of the Moon (1973)
- Public Enemy - It Takes a Nation of Millions to Hold Us Back (1965)
- Sex Pistols - Never Mind the Bollocks, Here's the Sex Pistols (1977)
- Stevie Wonder - Innervisions (1973)
- The Band - Music from Big Pink (1968)
- The Beach Boys - Pet Sounds (1966)
- The Beatles - Please Please Me (1963)
- The Beatles - Sgt. Pepper's Lonely Hearts Club Band (1967)
- The Jimi Hendrix Experience - Are You Experienced? (1967)
- The Rolling Stones - Exile on Main St. (1972)
- The Who - Who's Next (1971)
- U2 - The Joshua Tree (1987)

http://sobigdata.borgheresi.it/music_in_numbers/
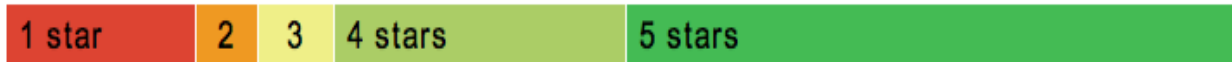
# Google Product Search

**HP Officejet 6500A Plus e-All-in-One Color Ink-jet - Fax / copier / printer / scanner**
**$89** online, **$100** nearby ★★★★☆ 377 reviews
September 2010 - Printer - HP - Inkjet - Office - Copier - Color - Scanner - Fax - 250 sh

## Reviews

**Summary** - Based on 377 reviews

| 1 star | 2 | 3 | 4 stars | 5 stars |
|--------|---|---|---------|---------|

What people are saying

| | | |
|---|---|---|
| ease of use | | "This was very easy to setup to four computers." |
| value | | "Appreciate good quality at a fair price." |
| setup | | "Overall pretty easy setup." |
| customer service | | "I DO like honest tech support people." |
| size | | "Pretty Paper weight." |
| mode | | "Photos were fair on the high quality mode." |
| colors | | "Full color prints came out with great quality." |

9

# Bing Shopping

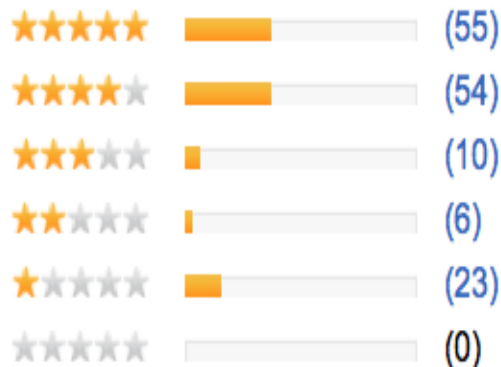## HP Officejet 6500A E710N Multifunction Printer

Product summary  Find best price  **Customer reviews**  Specifications  Related items

$121.53 - $242.39 (14 stores)

☐ Compare

Average rating ★★★☆☆ (144)

| | | |
|---|---|---|
| ★★★★★ | | (55) |
| ★★★★☆ | | (54) |
| ★★★☆☆ | | (10) |
| ★★☆☆☆ | | (6) |
| ★☆☆☆☆ | | (23) |
| ☆☆☆☆☆ | | (0) |

Most mentioned

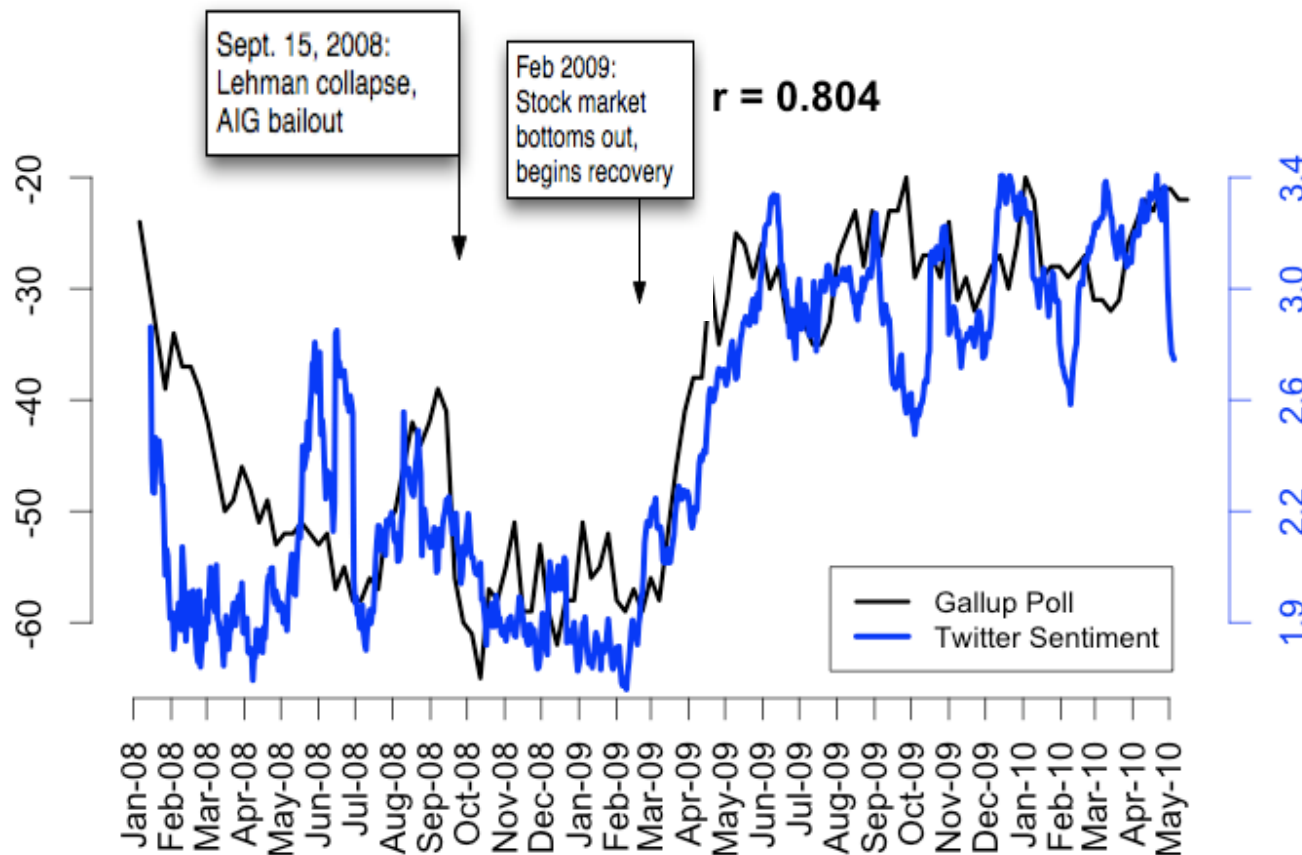| | | |
|---|---|---|
| Performance | | (57) |
| Ease of Use | | (43) |
| Print Speed | | (39) |
| Connectivity | | (31) |

More ▼

Show reviews by source

Best Buy (140)
CNET (5)
Amazon.com (3)

# Twitter sentiment versus Gallup Poll of Consumer Confidence

Brendan O'Connor, Ramnath Balasubramanyan, Bryan R. Routledge, and Noah A. Smith. 2010. From Tweets to Polls: Linking Text Sentiment to Public Opinion Time Series. In ICWSM-2010
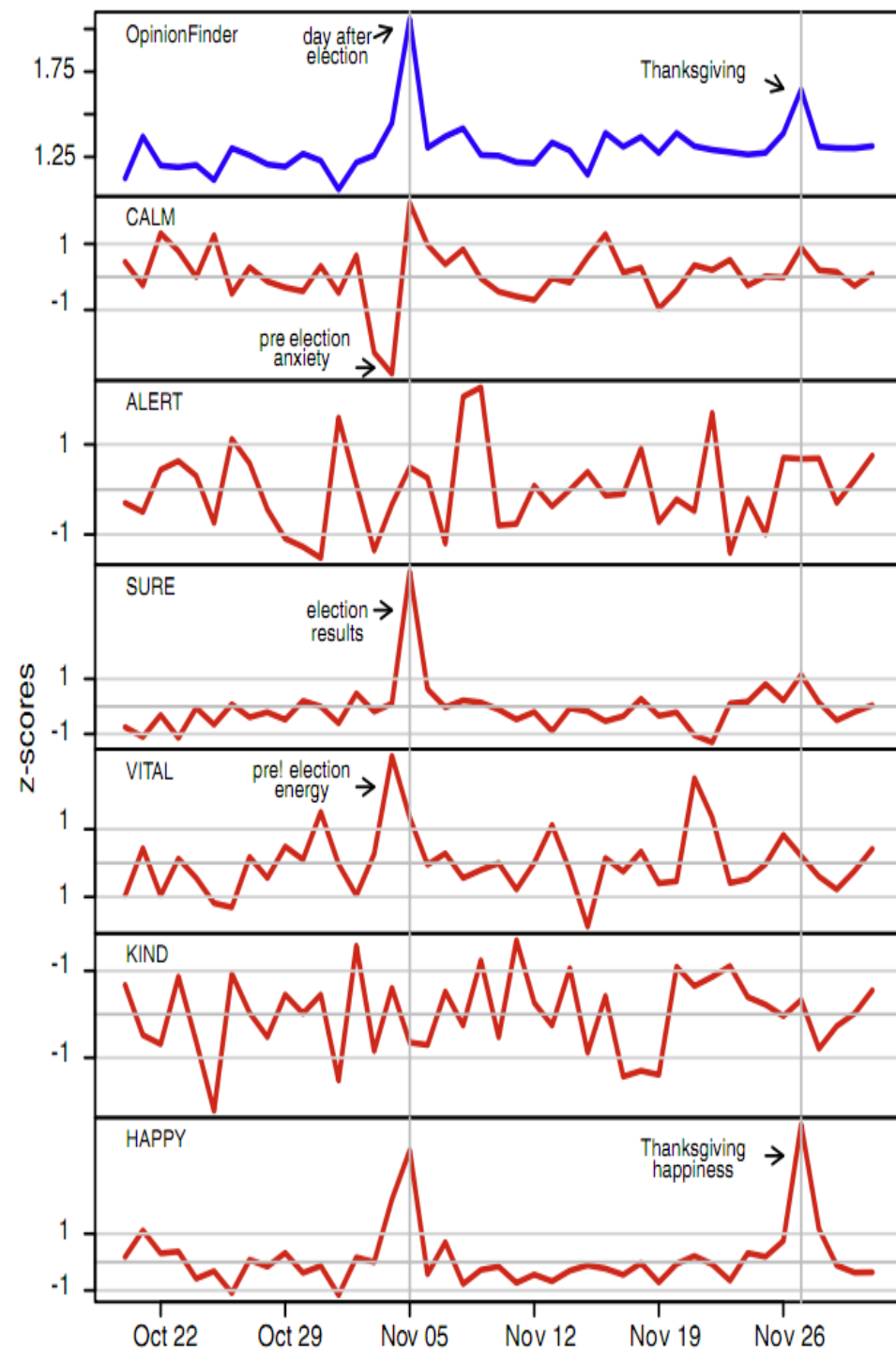
# Twitter sentir

Johan Bollen, Huina Mao, Xiaojun Zeng. 2011.
Twitter mood predicts the stock market,

Journal of Computational Science 2:1, 1-8.
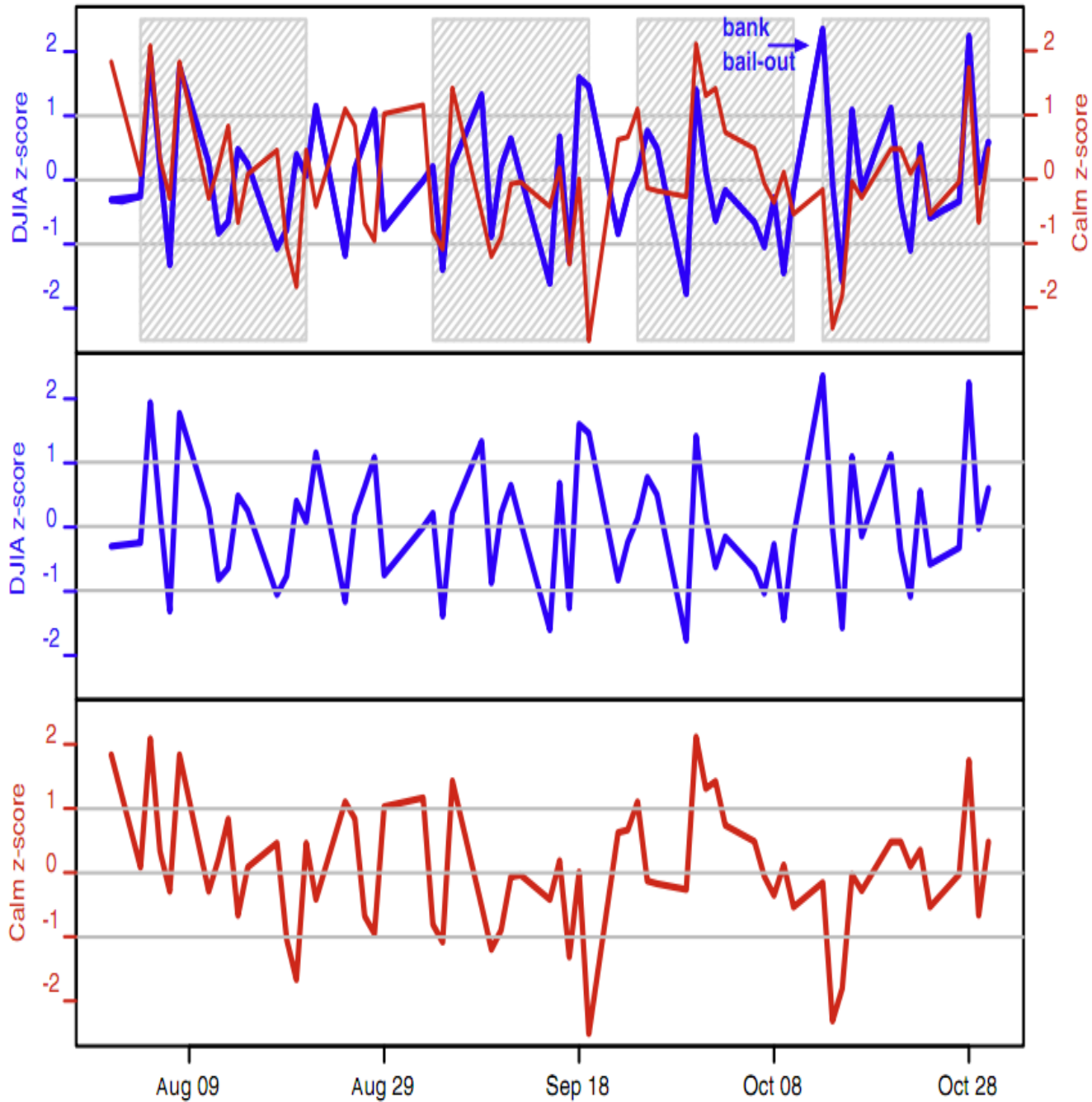10.1016/j.jocs.2010.12.007.

Opinion Finder ed GPOMOS, 7 Time series

Bollen et al. (2011)

CALM predicts DJIA 3 days later

At least one current hedge fund uses this algorithm

# Target Sentiment on Twitter

[Twitter Sentiment App](#)

Alec Go, Richa Bhayani, Lei Huang. 2009. Twitter Sentiment Classification using Distant Supervision

Type in a word and we'll highlight the good and the bad

"united airlines"    Search   Save this search

**Sentiment analysis for "united airlines"**

Sentiment by Percent | Sentiment by Count

Negative (68%)
Positive (32%)

Positive (11)
Negative (23)

jljacobson: OMG... Could @**United airlines** have worse customer service? W8g now 15 minut
Posted 2 hours ago

12345clumsy6789: I hate **United Airlines** Ceiling!!! Fukn impossible to get my conduit in this d
Posted 2 hours ago

EMLandPRGbelgiu: EML/PRG fly with Q8 **united airlines** and 24seven to an exotic destinatio
Posted 2 hours ago

CountAdam: FANTASTIC customer service from **United Airlines** at XNA today. Is tweet more
Posted 4 hours ago

# *Nowcasting: Predicting the present*

# The traps of big data

*Google Flu Trends* : search data can help predict the incidence of influenza-like diseases

Close relationship between number of people searching for flu-related topics and how many people have symptoms

Prediction models compared to real-world cases of flu

Hyunyoung Choi and Hal Varian. *Predicting the present with google trends*. Technical Report, 2009.

# How it works

- a time series is computed for about 50 million common queries entered weekly from 2003 to 2008.

- georeferenced by identifying the IP address associated with each search, the state in which this query was entered can be determined.

- linear model is used to compute the log-odds of Influenza-like illness (ILI) physician visit (official data) and the log-odds of ILI-related search query:

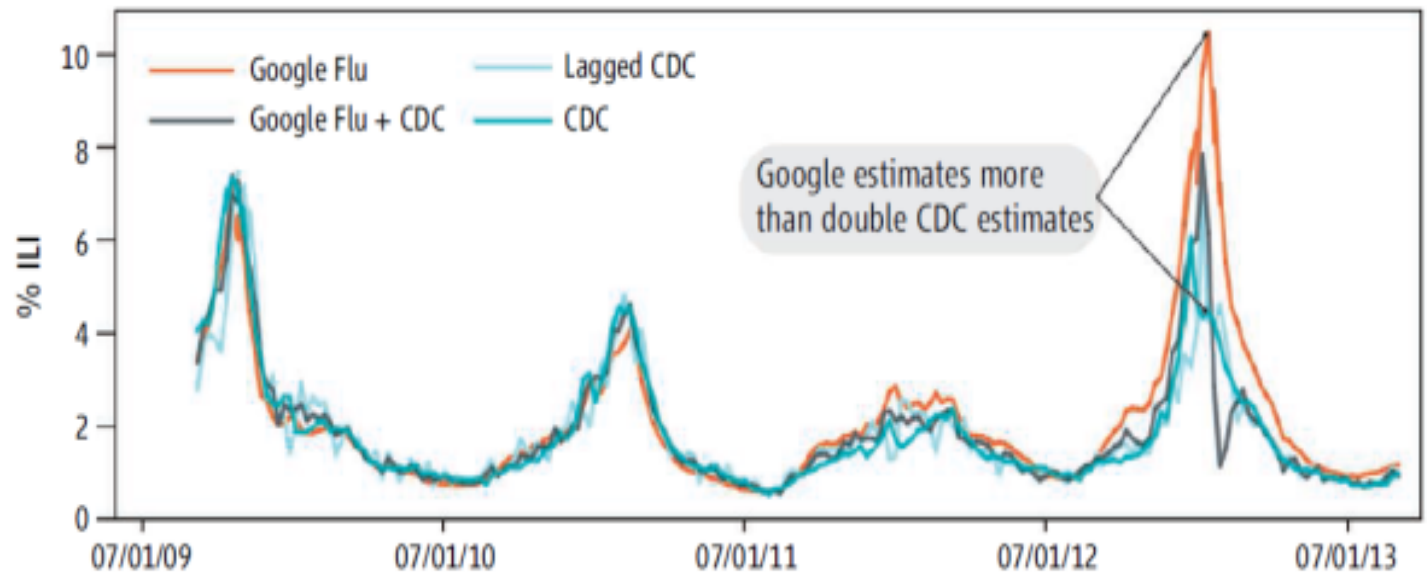$$\text{logit}(P) = \beta_0 + \beta_1 \times \text{logit}(Q) + \epsilon$$

# How it works

- Each of the 50 million queries is tested as Q to see if the result computed from a single query could match the actual history ILI data obtained from the U.S. Centers for Disease Control and Prevention (CDC).

- This process produces a list of top queries which gives the most accurate predictions of CDC ILI data when using the linear model. Then the top 45 queries are chosen because, when aggregated together, these queries fit the history data the most accurately.

- Finally, the trained model is used to predict flu outbreak across all regions in the United States.

# But..



- In 2009, completely missed the    non-seasonal influenza A-H1N1
- In 2013, double doctor visits than CDC

David Lazer, Ryan Kennedy, Gary King, and Alessandro Vespignani.
*The parable of google flu: Traps in big data analysis*. Science Magazine, 343(6176):1203–1205, 2014.

# Happiness as Subjective Well-Being Indicator

MEASURING HAPPYNESS

# Subjective Well-Being

- Perceptions and evaluations affect the way people face life and take advantage of opportunities in different ways.

- Subjective indicators are useful complement to strictly objective indicators, because they allow evaluating the possible differences between what people report and what it is captured by objective indicators.
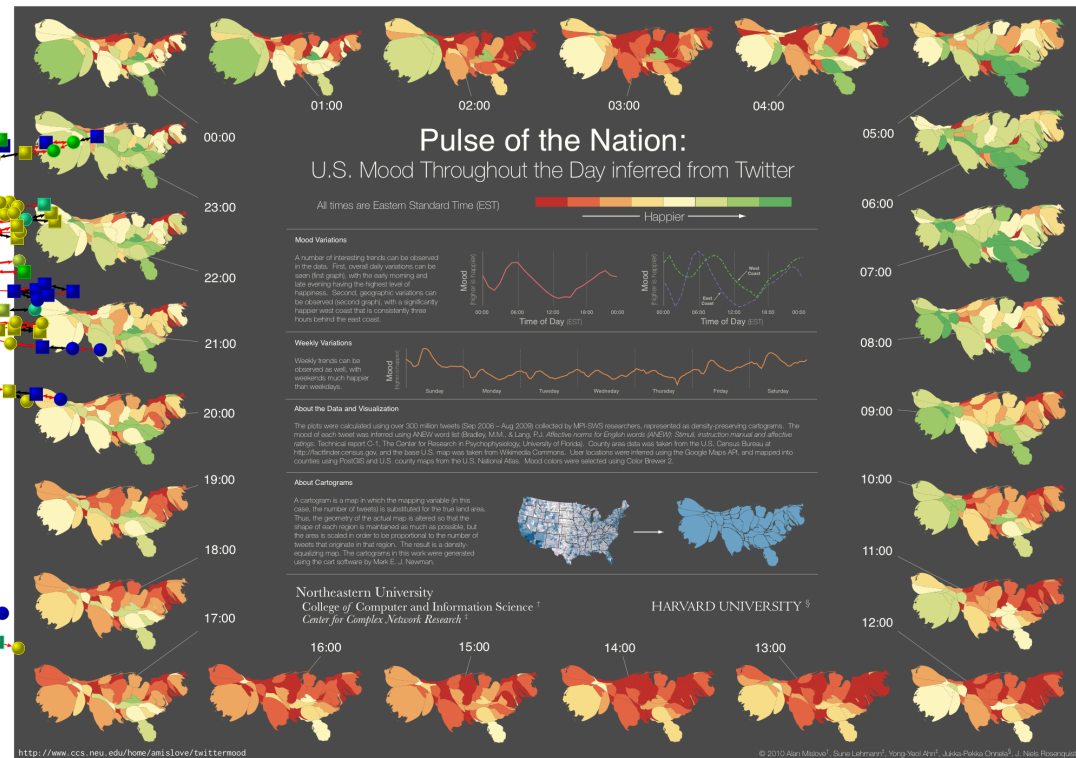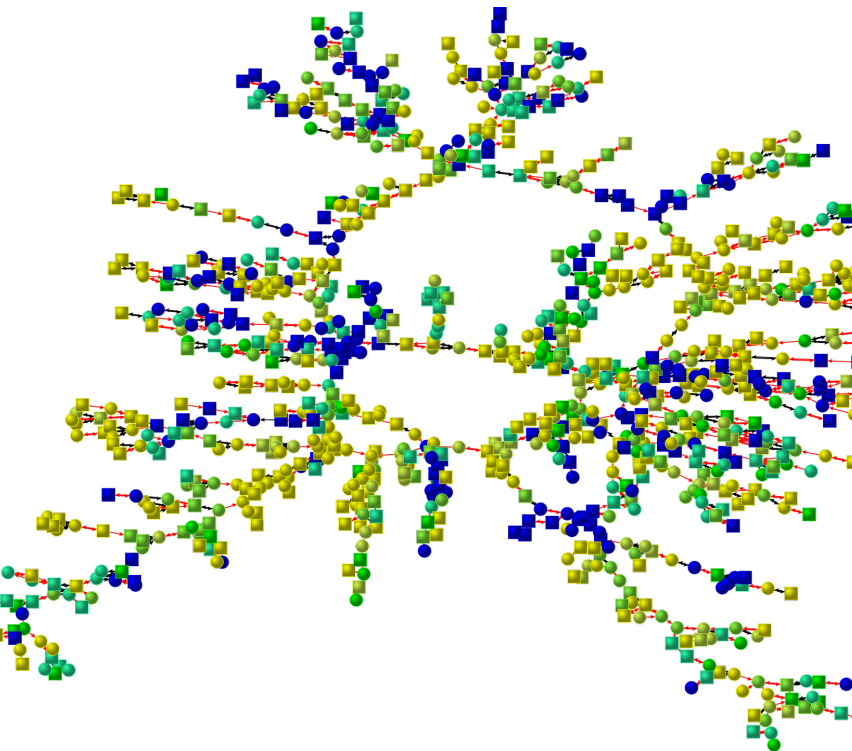
# Happiness as Subjective Well-Being Indicator

All the subjective well-being indicators observed in the literature come from surveys, and are a bag of ratings of questions like: life satisfaction, leisure time satisfaction, utilities closeness, etc.

From a data-driven point of view, the best way to estimate the level of subjective well-being is to estimate the level of happiness.

Happiness is intrinsically correlated with, and it is a consequence of, the percepted level of wellbeing driven by social relationships, health, work condition, etc.

Measuring the Happiness of Large-Scale Written Expression: Songs, Blogs, and Presidents.Peter Sheridan Dodds et al. Journal of Happiness Studies. 2010.
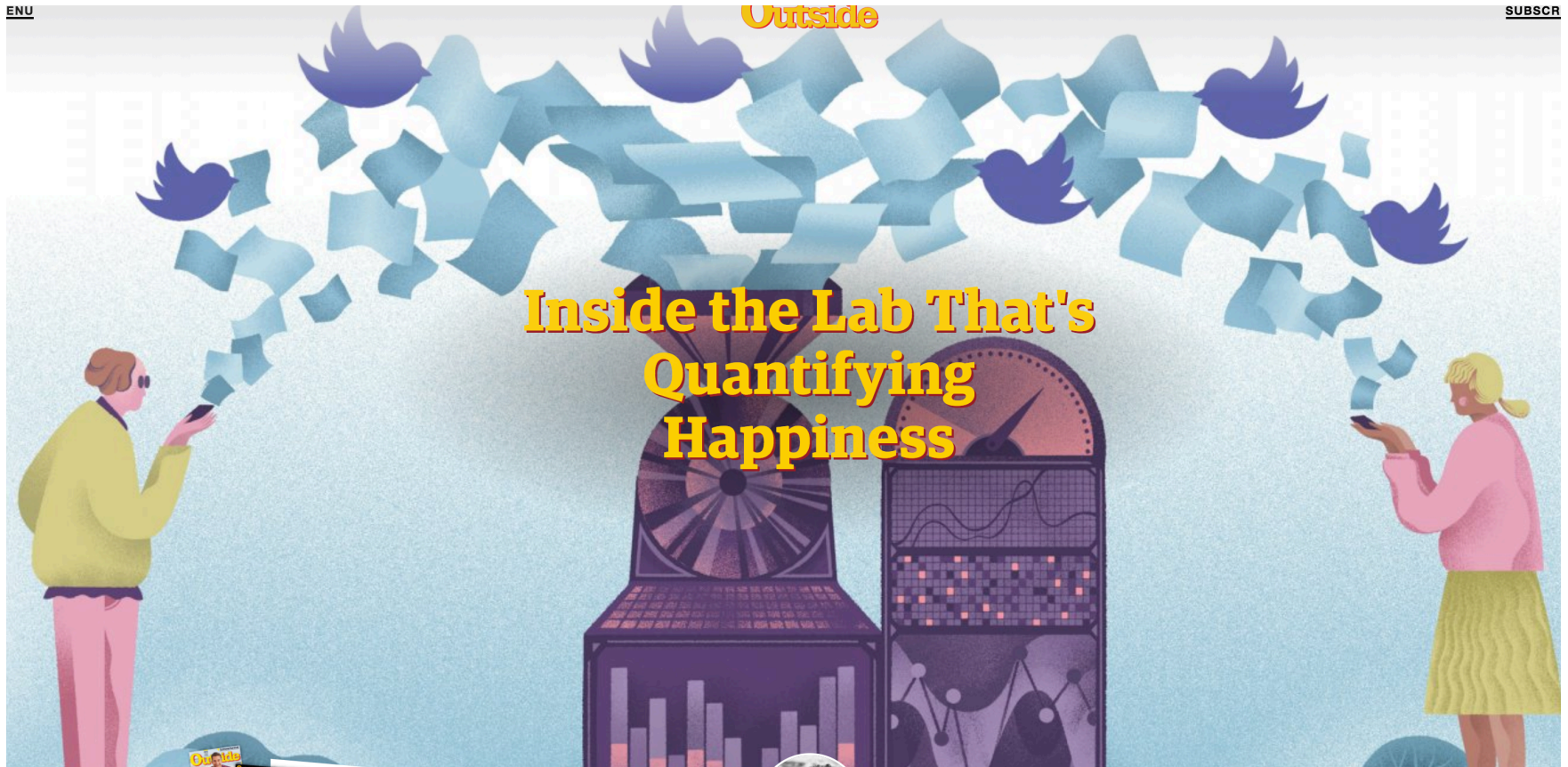
# Measuring Happyness

# Text Data Sources

- Every written text (books, songs, blogs, posts, etc.) reveals a positive or negative sentiment with respect to their content. Online Social Networks like Twitter and Facebook are the expression of our personal moods and our opinions.

- They are a very powerful data source to reveal the level of happiness in a certain region and in a certain *period*.

- Indeed one of the biggest issue of traditional indicators based on surveys is that they "measure" a certain variable in the instant the user answer the survey.

- The mood of a user and consequently the answer to a certain question could be different just a week after because of family or work problems for example.

**Inside the Lab That's Quantifying Happiness**

Measuring the Happiness of Large-Scale Written Expression: Songs, Blogs, and Presidents, Peter Sheridan Dodds Æ Christopher M. Danforth, J Happiness Stud (2010)

The Geography of Happiness: Connecting Twitter Sentiment and Expression, Demographics, and Objective Characteristics of Place Lewis Mitchell1 , Morgan R. Frank, Kameron Decker Harris1,2, Pete Sheridan Dodds, Christopher M. Danforth1

https://www.outsideonline.com/2230891/inside-lab-thats-quantifying-happiness

# A data-driven Happiness Indicator

ANEW: Affective Norms for English Words are 1034 words rated between 1-9 (good-bad) in a sort of happiness scale.

These words can be used to evaluate the level of happiness of a text called **valence**:

$$v_{\text{text}} = \frac{\sum_{i=1}^{n} v_i f_i}{\sum_{i=1}^{n} f_i}$$

# Valence of lyrics - example

## Lyrics for Michael Jackson's Billie Jean

"She was more like a beauty queen from a movie scene.

⋮

And mother always told me, be careful who you love. And be careful of what you do 'cause the lie becomes the truth.

Billie Jean is not my lover, She's just a girl who claims that I am the one.

⋮

| ANEW words | $v_k$ | $f_k$ |
|---|---|---|
| $k=1.$ love | 8.72 | 1 |
| 2. mother | 8.39 | 1 |
| 3. baby | 8.22 | 3 |
| 4. beauty | 7.82 | 1 |
| 5. truth | 7.80 | 1 |
| 6. people | 7.33 | 2 |
| 7. strong | 7.11 | 1 |
| 8. young | 6.89 | 2 |
| 9. girl | 6.87 | 4 |
| 10. movie | 6.86 | 1 |
| 11. perfume | 6.76 | 1 |
| 12. queen | 6.44 | 1 |
| 13. name | 5.55 | 1 |
| 14. lie | 2.79 | 1 |

$$v_{\text{text}} = \frac{\sum_k v_k f_k}{\sum_k f_k}$$
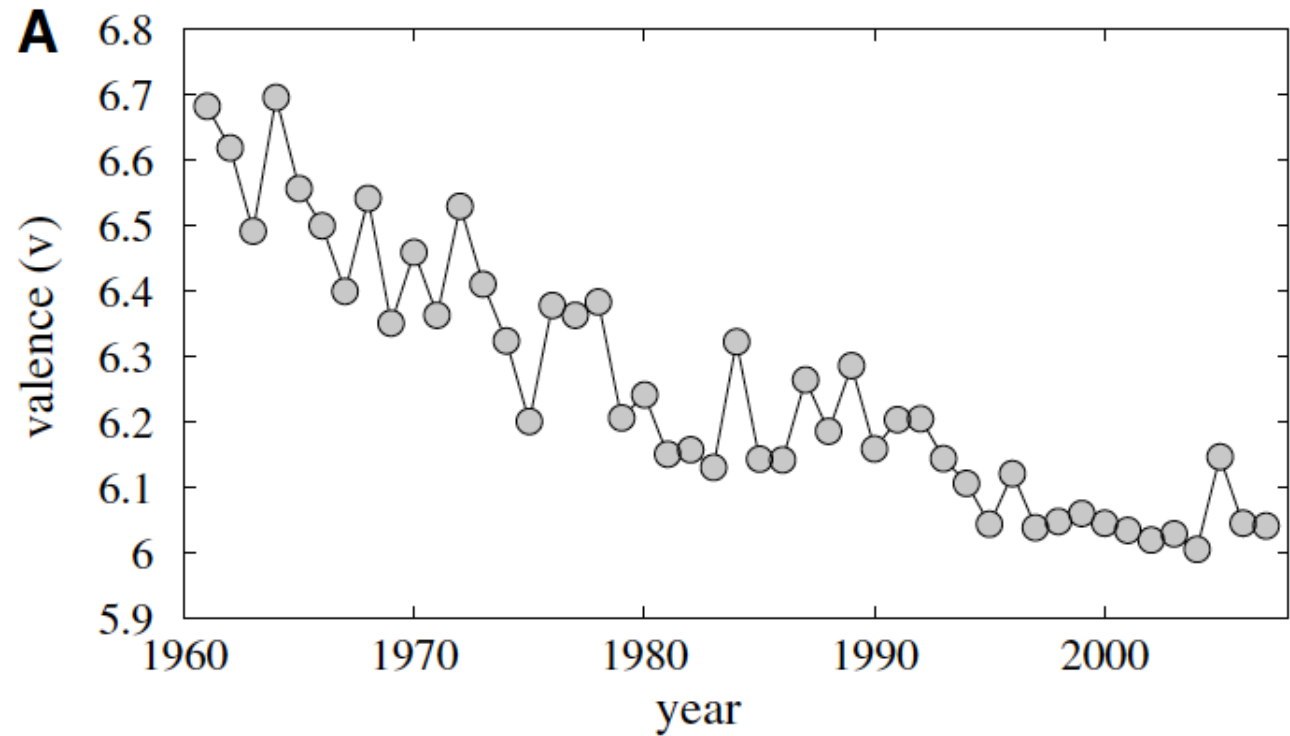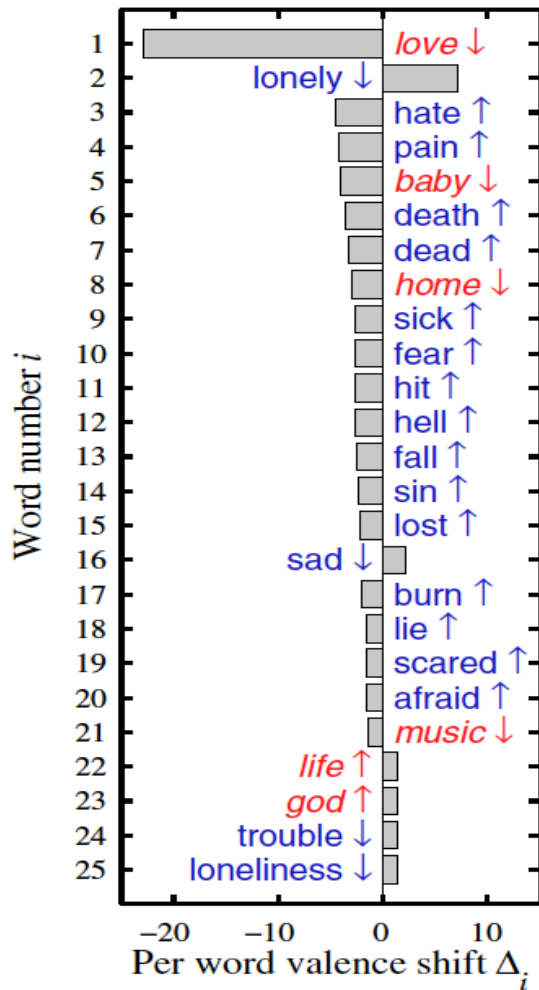
$v_{\text{Billie Jean}} = 7.1$
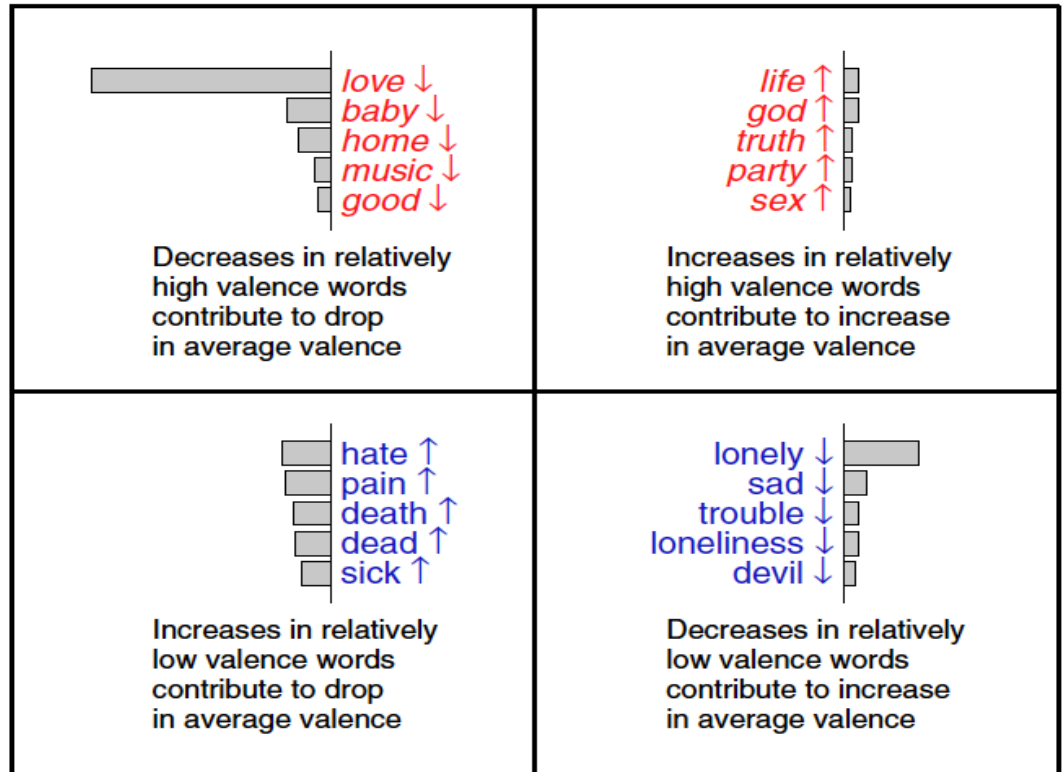
$v_{\text{Thriller}} = 6.3$

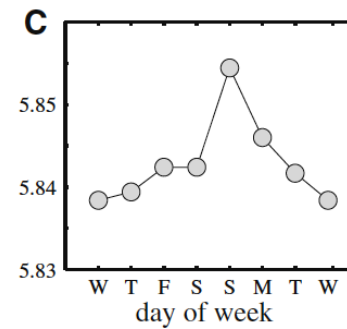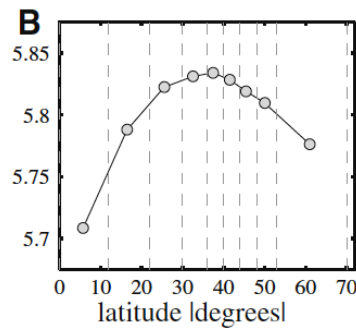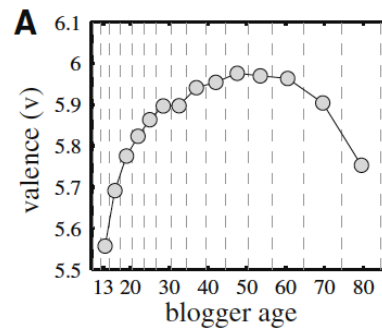$v_{\text{Michael Jackson}} = 6.4$

# Song lyrics valence across decades

Key:

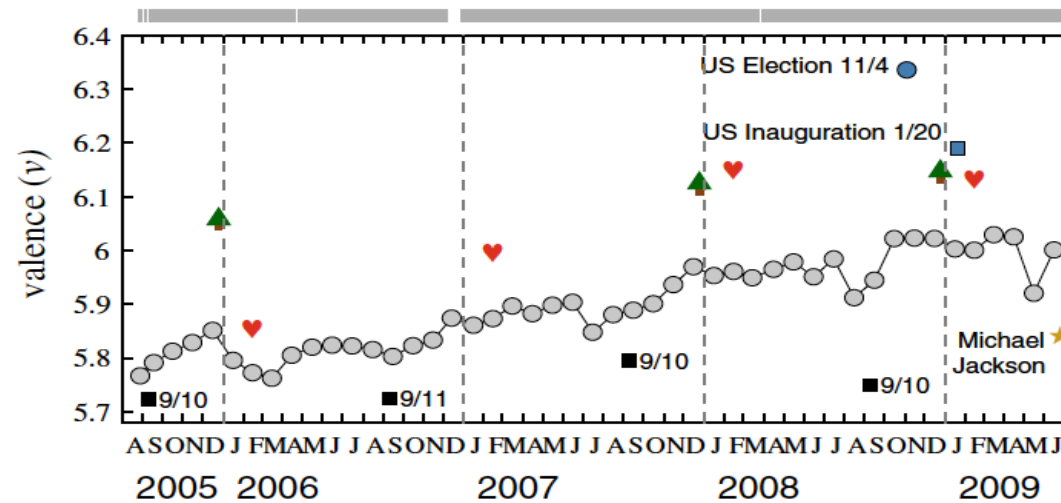Decreases in relatively high valence words contribute to drop in average valence

Increases in relatively high valence words contribute to increase in average valence

Increases in relatively low valence words contribute to drop in average valence
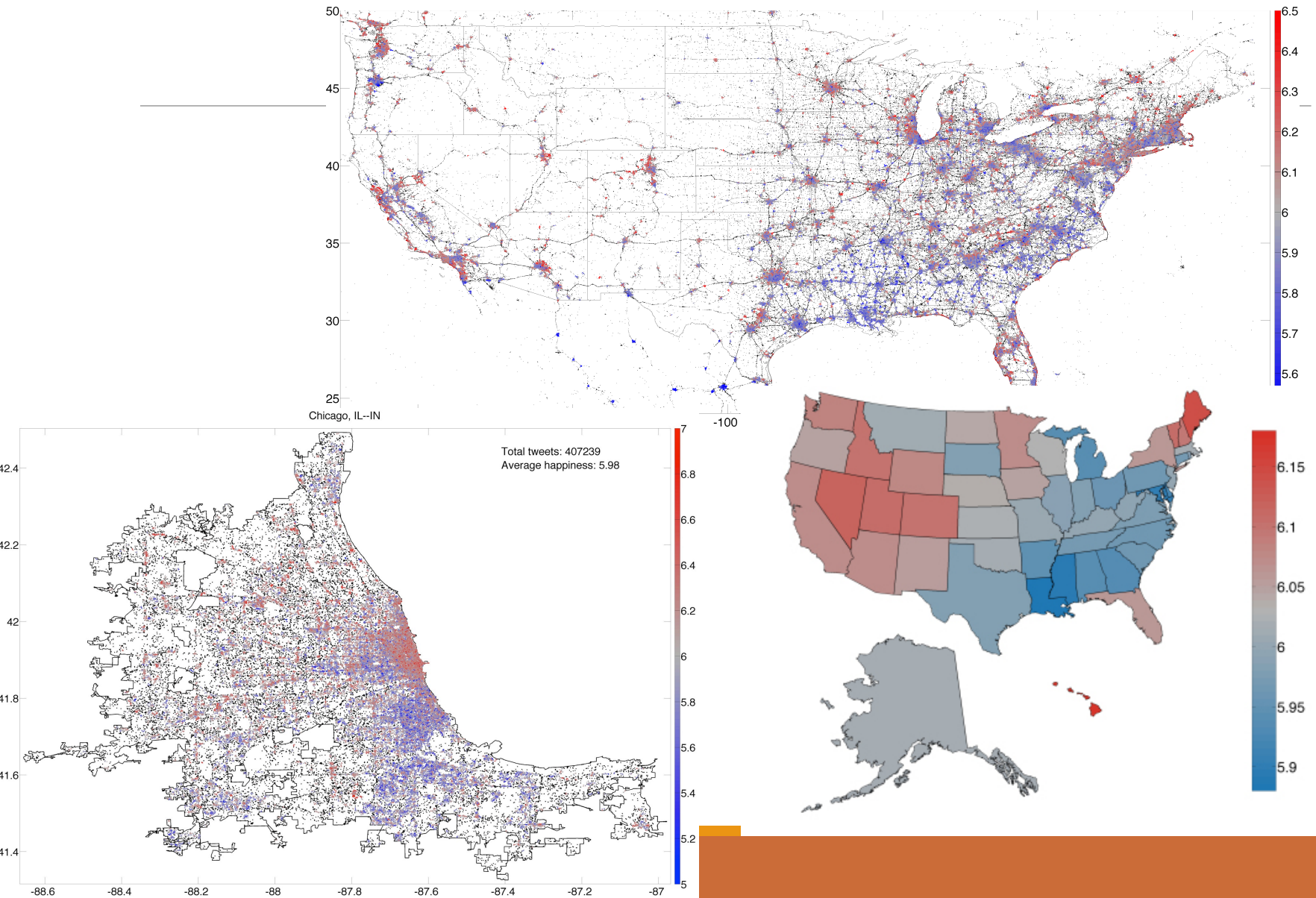
Decreases in relatively low valence words contribute to increase in average valence

# Blog valence

# Twitter hedonometer

# Average Happiness for Twitter

# PARTIAMO DALLE PAROLE: DOODS

| Rapimento | 2.76 |
|-----------|------|
| Dolore | 2.46 |
| Angosciato | 2.12 |
| Sanguinoso | 2.90 |

| Spiaggia | 8.03 |
|----------|------|
| Amore | 8.72 |
| Bacio | 8.26 |
| Casa | 7.91 |

# COME CALCOLIAMO IL MOOD?

- Frase: "io amo le spiagge"
- Lemmatizzazione:

  io → io

  amo → amare

  le → la

  spiagge → spiaggia
- Applicazione della formula di Dodds

# LEMMATIZZAZIONE:
## RIPORTARE LE FORME FLESSE A LEMMI

Io amo le spiagge → io amare la spiaggia

Spiaggia     8.03

Io amare la spiaggia

Amore   8.72

8.72          8.03

# DOODS:

**CALCOLIAMO LA SOMMATORIA DEGLI SCORE TROVATI MOLTIPLICATI PER LA LORO FREQUENZA DIVISO LA SOMMATORIA DELLE FREQUENZE DI TUTTI I LEMMI TROVATI**

$$\frac{\sum_{i=1}^{n} v_i f_i}{\sum_{i=1}^{n} f_i} = \frac{8.72 * 1 + 8.03 * 1}{1 + 1} = 8.3$$

# Measuring the "salad bowl" -Superdiversity on Twitter-

## ALINA SÎRBU

WITH

LAURA POLLACCI, FOSCA GIANNOTTI, DINO PEDRESCHI

*UNIVERSITY OF PISA*

*ISTI CNR*

# Outline

Superdiversity - a novel index based on sentiment

Validation against immigration data

Comparison with other indices

Discussion: is nowcasting possible?

# Superdiversity

**Superdiversity** - a new level of cultural diversity due to immigration and cultural differences among immigrants themselves (Vertovec, 2007)

Measuring superdiversity - difficult task

◦ Cultural diversity - number of languages spoken in a region

◦ Immigration rates - official statistics - generally low time and space resolution

◦ Use social big data

  ◦ Here: geolocalised tweets

# Superdiversity Index (SI)

Main idea:

◦ different cultures assign **different emotional valence** to different words.

◦ compute **SI** as a **distance** between the '**standard**' emotional valence of a set of words and the '**used**' valence in the population of a region

  ◦ standard - manually tagged lexicon - ANEW (Bradley and Lang 1999)

  ◦ used - estimate from Twitter data

# Estimating sentiment valences of words on Twitter

Using the algorithm from *Pollaci et al. 2017.*

Extend a sentiment-tagged lexicon using a Twitter corpus

◦ built to enhance lexicon-based sentiment analysis on Twitter

◦ starts from a small seed lexicon with sentiment valences

◦ builds a co-occurrence network of words from the Twitter corpus

◦ sentiment valences diffuse from the seed to the other words

Laura Pollacci, Alina Sîrbu, Fosca Giannotti, Dino Pedreschi, Claudio Lucchese,and Cristina Ioana Muntean. 2017. Sentiment Spreading: An Epidemic Model for Lexicon-Based Sentiment Analysis on Twitter. In Conference of the Italian Association for Artificial Intelligence. Springer, 114–127.
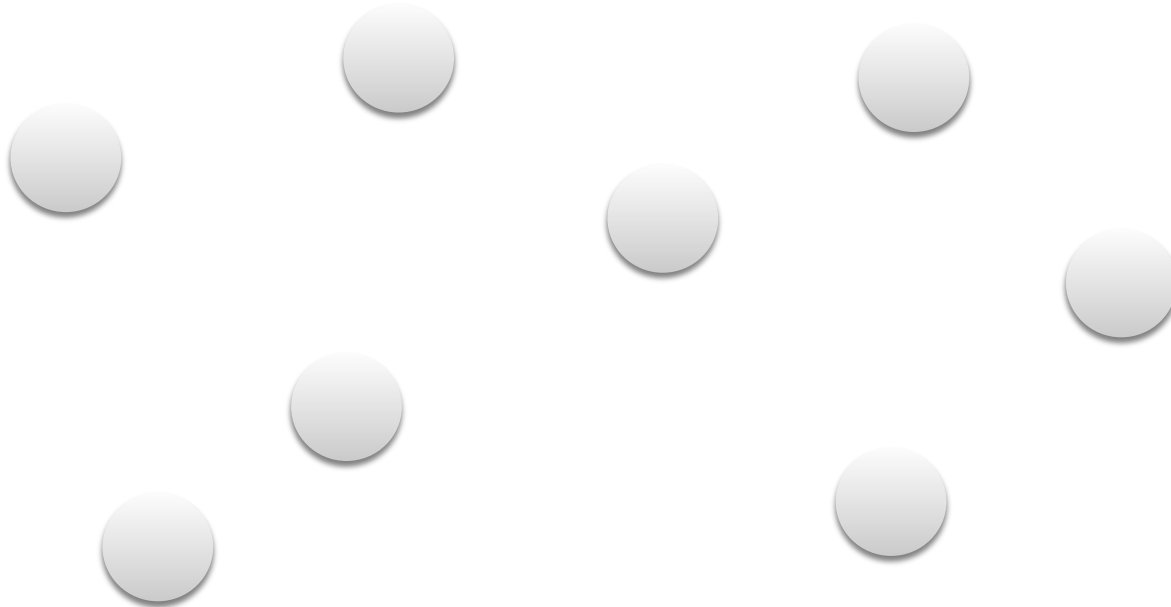
# Sentiment Spreading: An Epidemic Model for Lexicon-Based Sentiment Analysis

- Extend the dictionary and assign valence to tweet using a epidemic based approach (opinion dynamics like)
- Network of words - words that co-occur in a tweet are connected
- Some words have a valence, the other words take the valence of the neighbours (mean)
- After many iterations the system converges to a stable valence
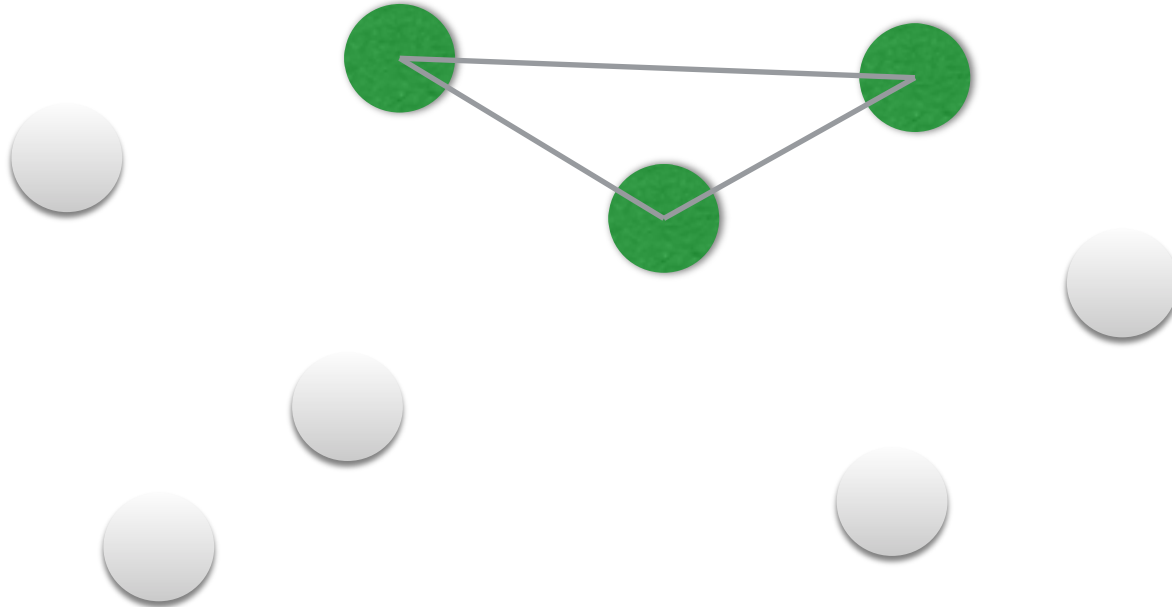- Method helps enlarge the dictionary, classify tweets but also characterise the group where the tweets came from
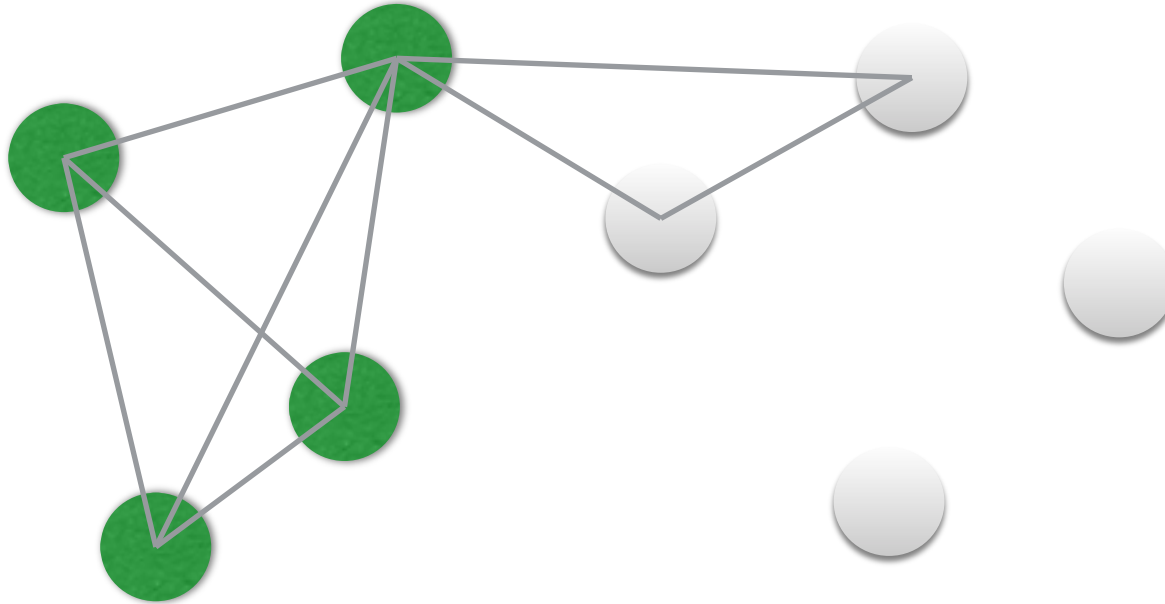
# Estimating sentiment valences of words on Twitter
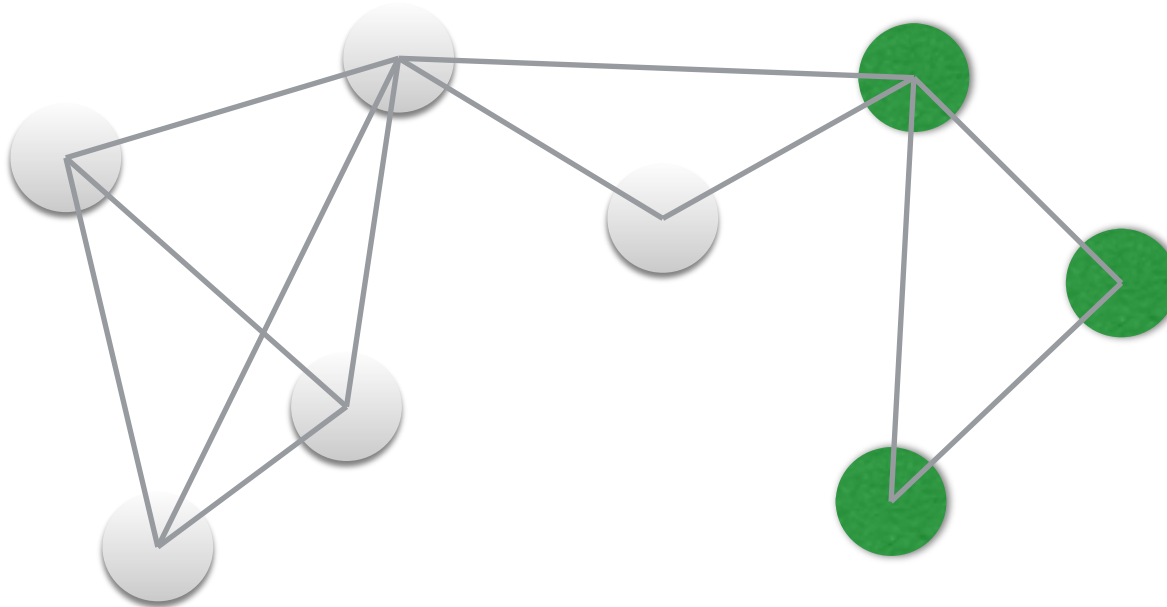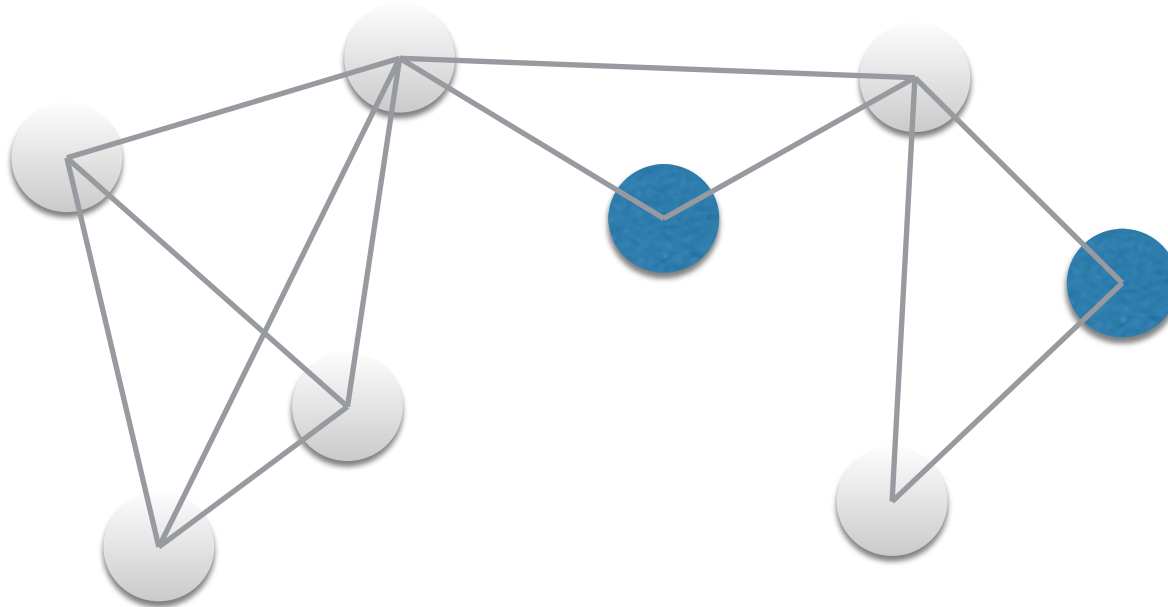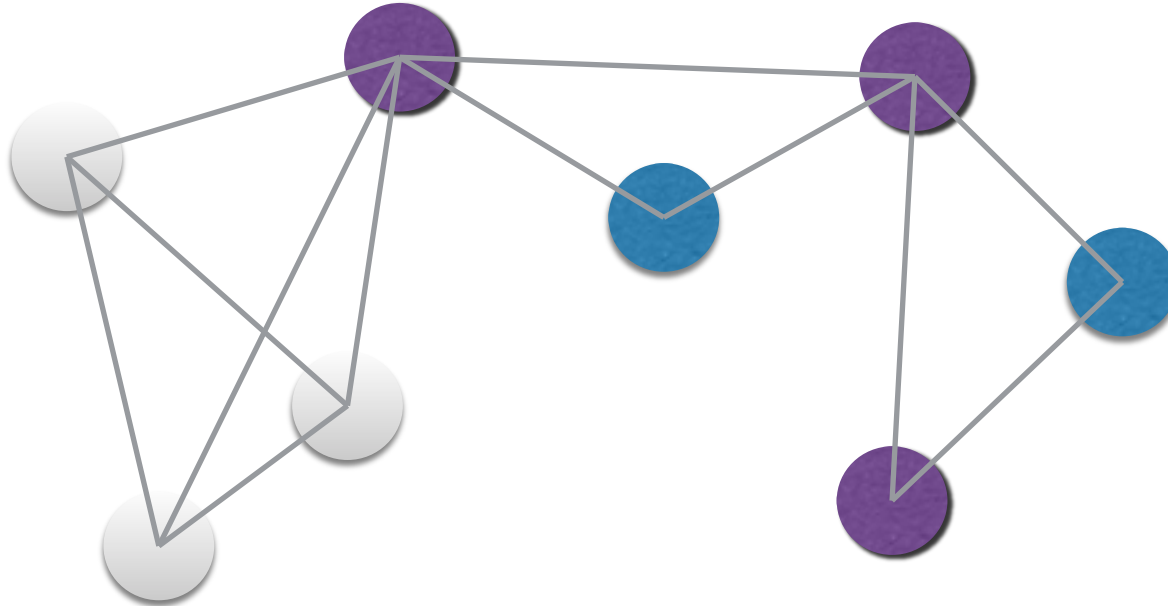
# Estimating sentiment valences of words on Twitter

# Estimating sentiment valences of words on Twitter

# Estimating sentiment valences of words on Twitter

# Estimating sentiment valences of words on Twitter

# Estimating sentiment valences of words on Twitter

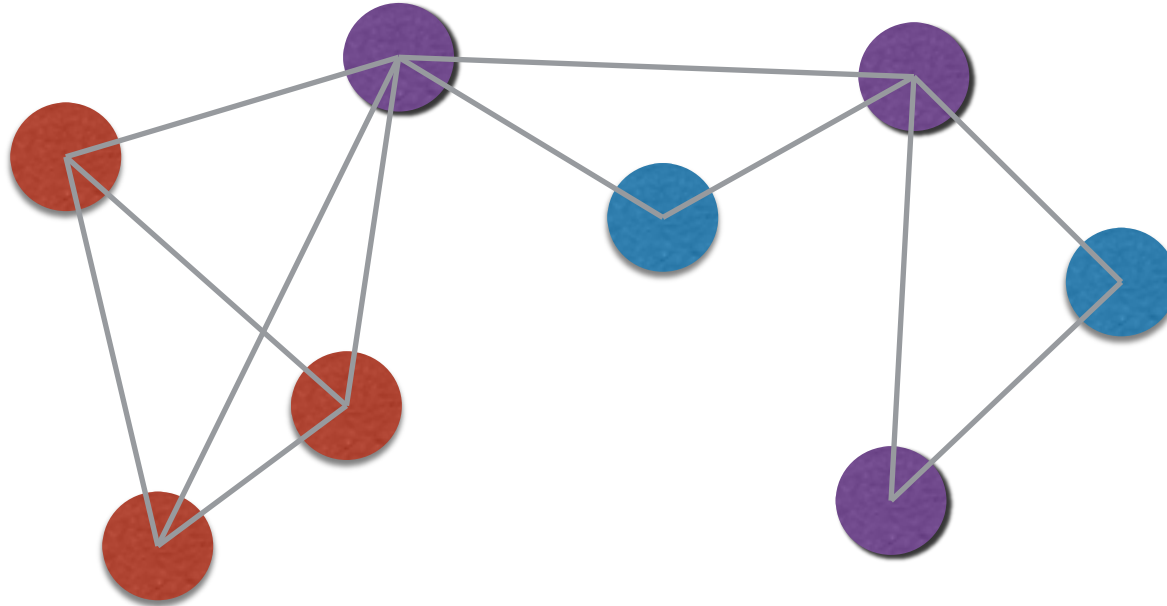# Estimating sentiment valences of words on Twitter

# Estimating sentiment valences of words on Twitter

The resulting dictionary

◦ is larger - enhanced sentiment analysis on Twitter

◦ depends on the way language is used - the new valences are population dependent

  ◦ we use the new valences as estimates for the real emotional content of the words in the population - **compute the distance to a manually tagged lexicon**
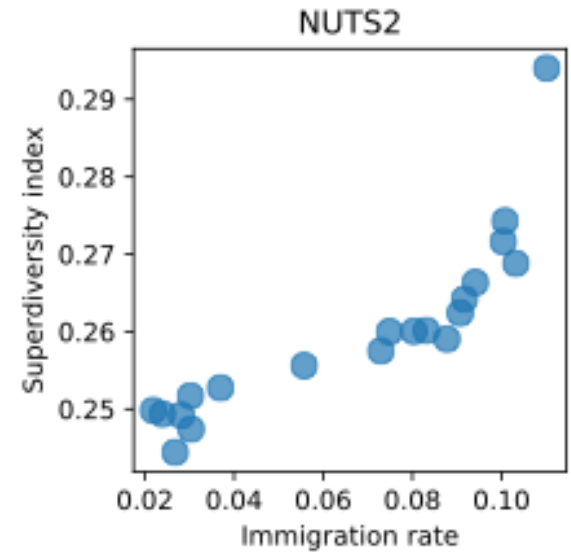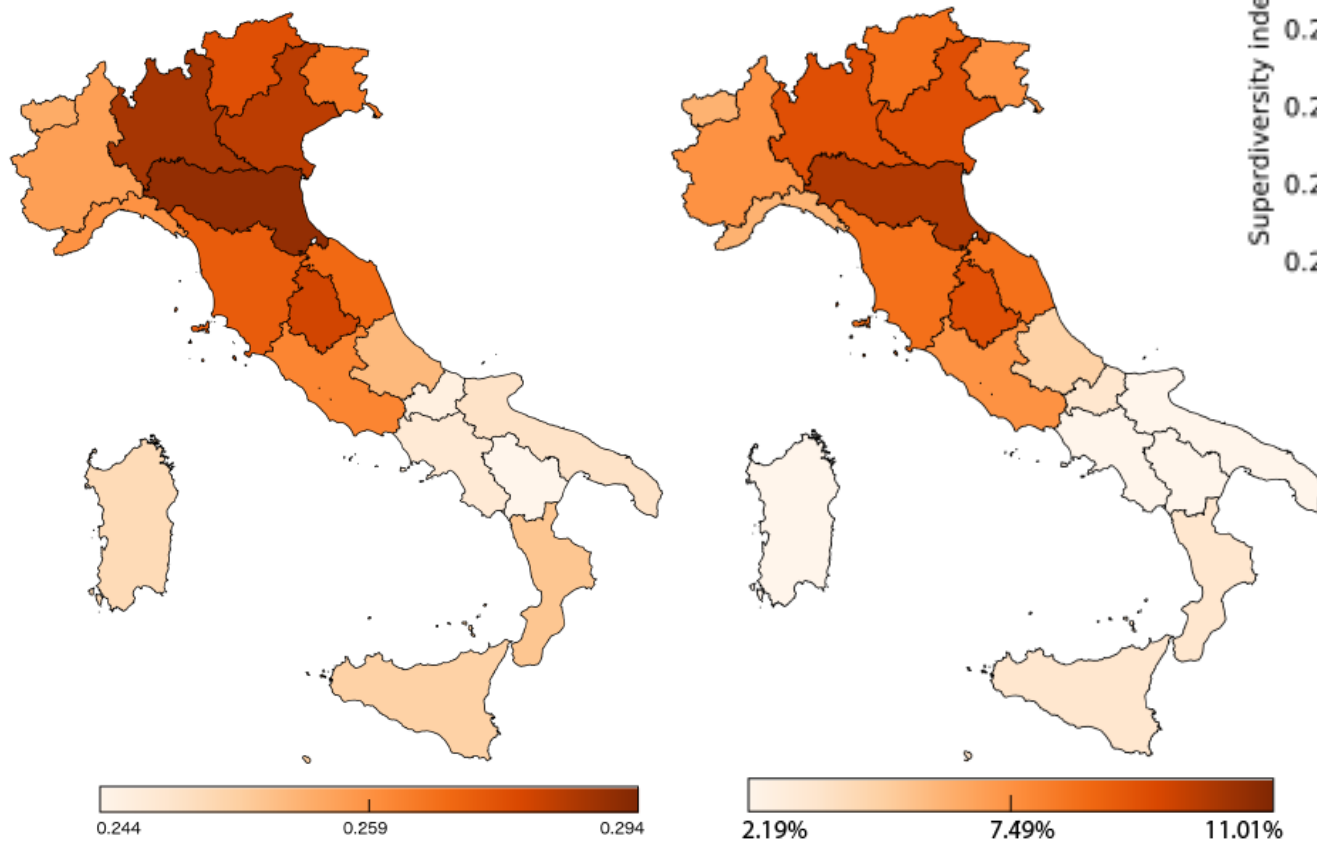
Compute SI for different geographical regions

◦ 10% of Geolocalised tweets for 3 months (August-October 2015)
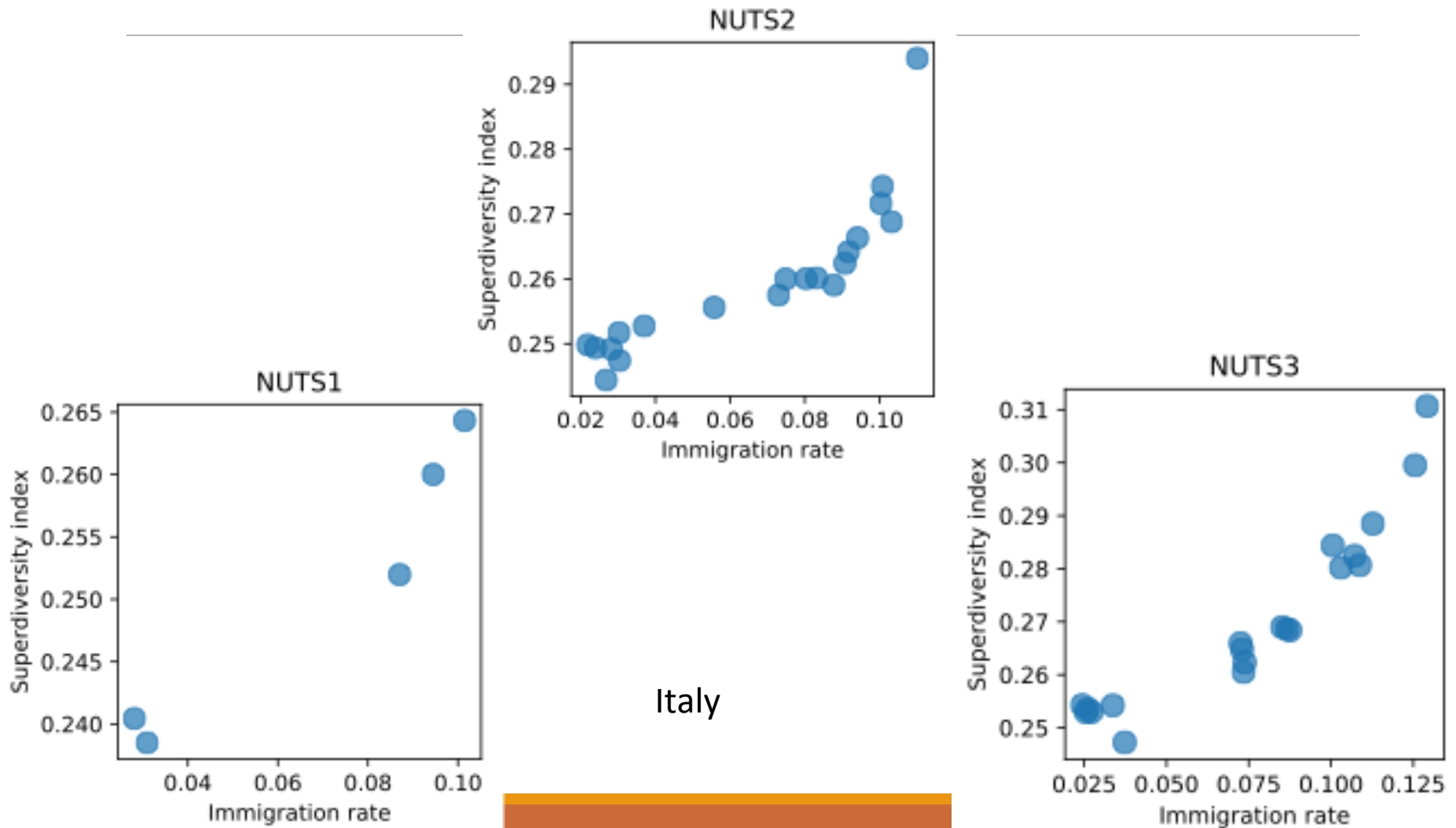
Compare SI with immigration rates

◦ JRC D4I dataset - immigration rates at various NUTS* levels  (https://bluehub.jrc.ec.europa.eu/datachallenge/)

◦ We analyse NUTS1, NUTS2, NUTS3 for Italy and UK

*Nomenclature of Territorial Units for Statistics

NUTS2

NUTS1

NUTS3

Italy

NUTS2

NUTS1

NUTS3
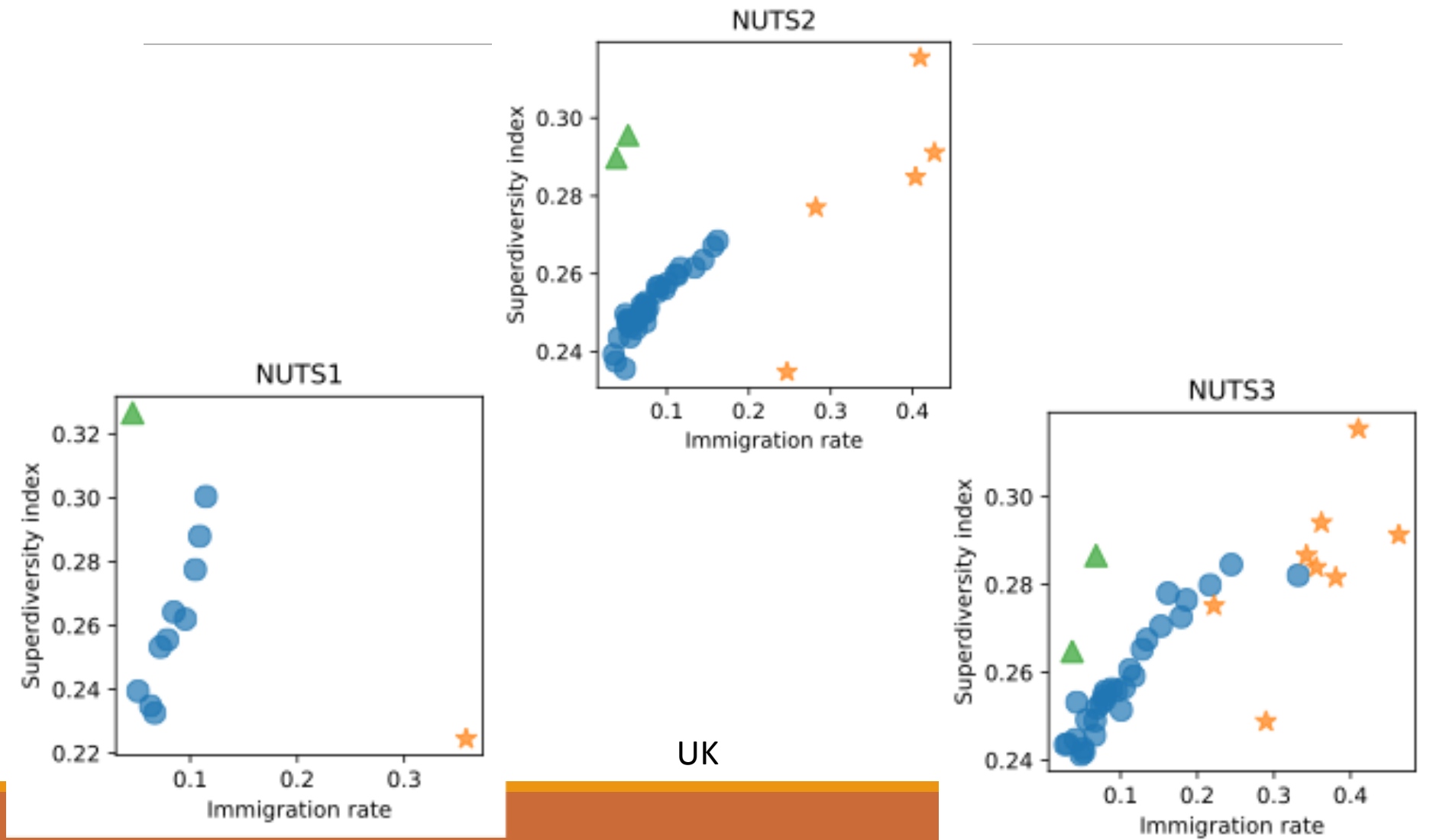
UK

Compare with other possible indices and null model
- ◦ Null model obtained by shuffling tweets between regions, but maintaining the number of tweets
- ◦ Other indices:
  - ◦ Number of tweets (/capita), **number of languages**, entropy of languages, type-token ratio (TTR)
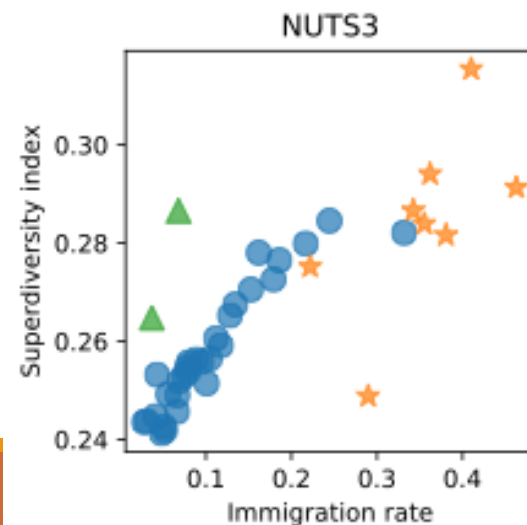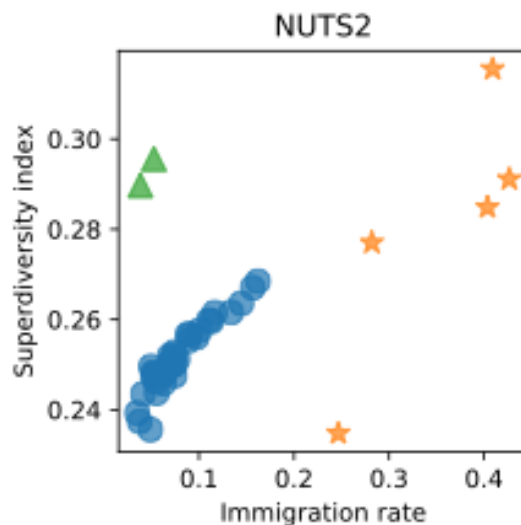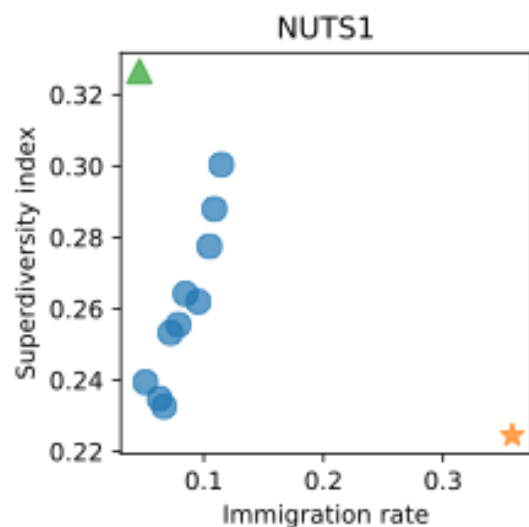
Italy

| Geographical level | SI | null model SI | number of Italian tweets | number of Italian tweets per capita | number of languages | language entropy | TTR |
|---|---|---|---|---|---|---|---|
| NUTS1 (5 regions) | **0.963** | -0.437 | 0.735 | 0.696 | 0.183 | -0.585 | -0.727 |
| NUTS2 (20 regions) | **0.859** | 0.143 | 0.279 | 0.282 | 0.304 | 0.099 | -0.243 |
| NUTS3 (20 regions) | **0.924** | 0.082 | 0.081 | -0.148 | 0.216 | 0.021 | 0.091 |

UK (except for London and Northeast England)

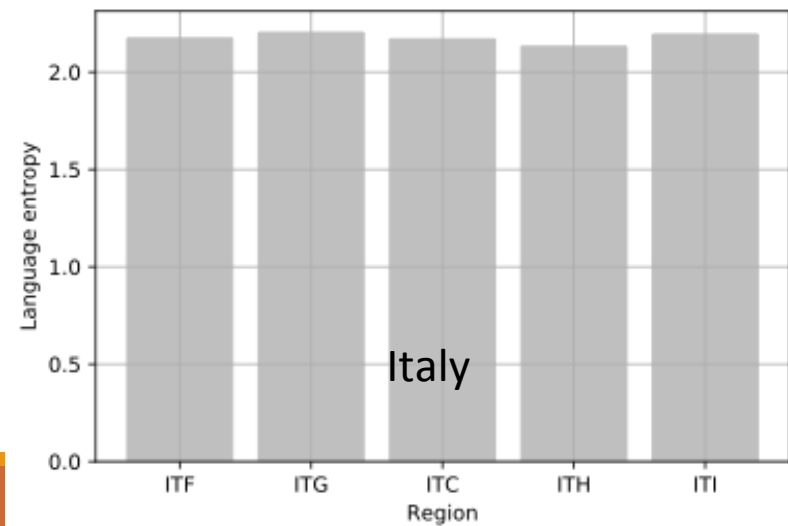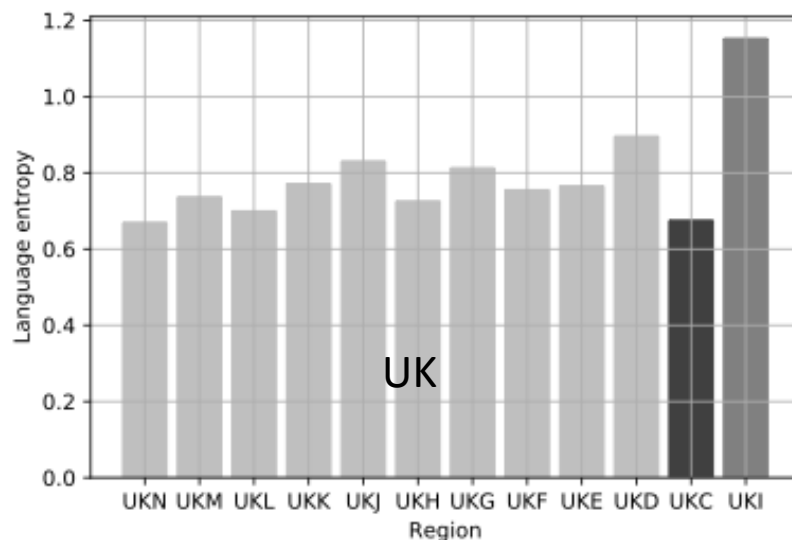| Geographical level | SI | null model SI | number of English tweets | number of English tweets per capita | number of languages | language entropy | TTR |
|---|---|---|---|---|---|---|---|
| NUTS1 (10 regions) | **0.943** | -0.236 | 0.328 | -0.520 | 0.519 | 0.481 | -0.005 |
| NUTS2 (40 regions) | **0.941** | -0.137 | 0.332 | 0.007 | 0.362 | 0.288 | -0.340 |
| NUTS3 (40 regions) | **0.928** | -0.221 | 0.141 | 0.049 | 0.322 | 0.529 | 0.147 |

What about London and Northeast England?

◦ SI appears to have different ranges in different region groups (and also spatial resolutions)

◦ Pre-clustering to identify these groups?

What about London and Northeast England?

◦ Pre-clustering to identify these groups?

  ◦ Language entropy - some information

  ◦ Need to identify further features - e.g. population density, local dialects

Can we use **superdiversity** to **nowcast immigration**?

◦ Correlations are, in general, very large.

◦ By adding pre-clustering and other correcting factors, within a ML model, it should be possible.

◦ We can estimate immigration levels at various resolutions where they are not available or up to date (clandestine immigration?)

Contact: alina.sirbu@unipi.it

References:

◦ Pollacci, Laura, Alina Sîrbu, Fosca Giannotti, Dino Pedreschi, Claudio Lucchese, and Cristina Ioana Muntean. *"Sentiment Spreading: An Epidemic Model for Lexicon-Based Sentiment Analysis on Twitter."* In AI*IA 2017.

◦ Pollacci, Laura, Alina Sîrbu, Fosca Giannotti, Dino Pedreschi, *"Measuring the 'Salad Bowl'- Superdiversity on Twitter."* Submitted.

VOICES
from the Blogs

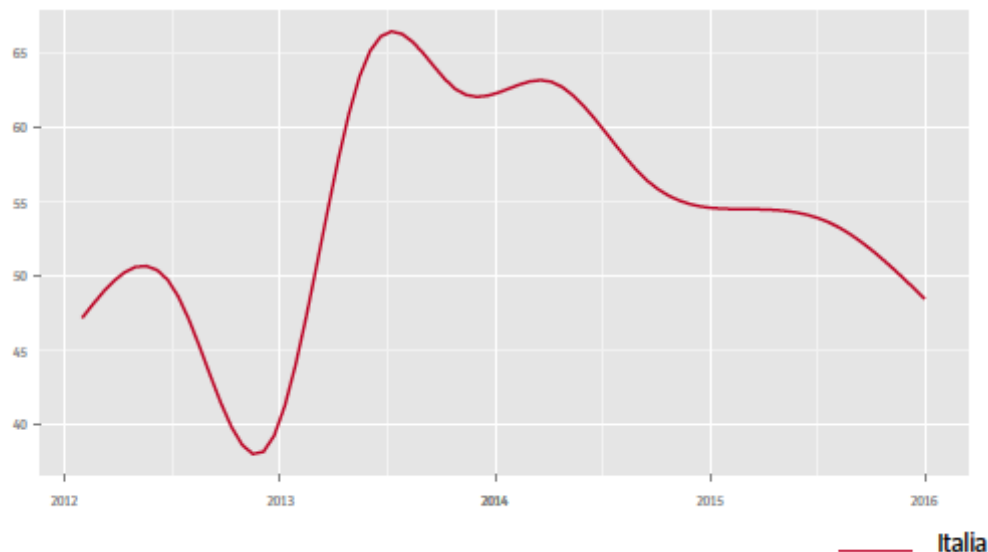# iHAPPY 2015

di
Andrea Ceron
Luigi Curini
Stefano M. Iacus

Politics and Big Data: Nowcasting and Forecasting Elections with Social Media. By Andrea Ceron, Luigi Curini, Stefano Maria
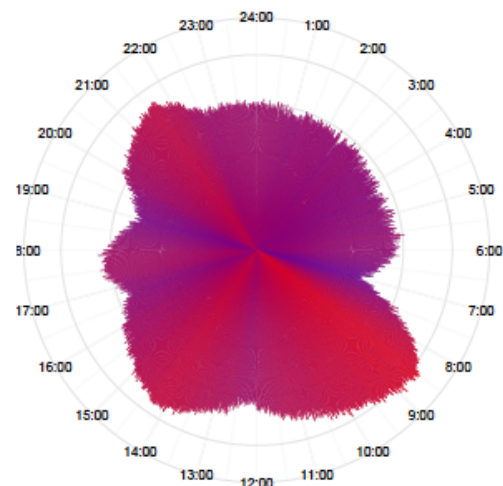
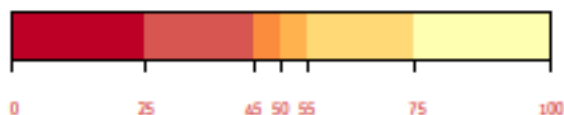# 2012-15 ANDAMENTO DI *IHAPPY* IN ITALIA

% di tweet felici



Italia

## L'ORA ITALIANA PIÙ FELICE: *IHAPPY* MINUTO PER MINUTO



# 2015 CALENDARIO DELLA TWITTER-FELICITÀ



| | GENNAIO | FEBBRAIO | MARZO | APRILE M | AGGIO | GIUGNO L | UGLIO | AGOSTO | SETTEMBRE O | TTOBRE N | OVEMBRE D | ICEMBRE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DOMENICA | | | | | | | | | | | | |
| LUNEDÌ | | | | | | | | | | | | |
| MARTEDÌ | | | | | | | | | | | | |
| MERCOLEDÌ | | | | | | | | | | | | |
| GIOVEDÌ | | | | | | | | | | | | |
| VENERDÌ | | | | | | | | | | | | |
| SABATO | | | | | | | | | | | | |

0    25    45 50 55    75    100

## LA CLASSIFICA REGIONALE

(La freccia indica se una regione è migliorata
o peggiorata nella classifica rispetto al 2014)

| Umbria | Puglia | Trentino Alto Adige | Toscana | Marche |
|---|---|---|---|---|
| 54,9% | 54,4% | 54% | 53,6% | 53,3% |
| Emilia Romagna | Piemonte | Liguria | Friuli V. G. | Sicilia |
| 53,3% | 53,2% | 53,1% | 53% | 52,1% |
| Lombardia | Sardegna | Calabria | Lazio | Veneto |
| 52% | 51,8% | 51,7% | 51,6% | 51,6% |
| Campania | Abruzzo | Basilicata | Valle d'Aosta | Molise |
| 51,5% | 51,2% | 50,7% | 48,6% | 47,6% |

# Text Analysis & Social Media

...dal rumore all'informazione

# Come analizzare i social?

**Re Tweet** @re_assoluto · 30 apr
Doniamo una mamma di Baltimora a ognuno dei ragazzi di #noexpo

| RETWEET | PREFERITI |
|---------|-----------|
| 380 | 387 |

05:17 - 30 apr 2015 · Dettagli

~~Mario Marco~~ @~~...~~ · 1 mag
I milanesi che, con spugna e sapone, ripuliscono i graffiti dei #noexpo sono il barlume di civiltà che ci serviva.

199    205    Foto

1) Non basta contare le mentions: Ma distinguere contenuto, ironia, giochi di parole, rumore, …

# Come analizzare i social?



## Semantic rules do work ?

◉ *Language evolves continuously: one cannot code all possible semantic rules unless reading the posts !!!*

???

"horses and bayonets" ?

ironic !

## Why human and not ontological dictionaries?

◉ *"What a nice rip-off"* ("che bella fregatura")

50% **positive** & 50% **negative**
=
misclassification

100% **negative**
=
no misclassification

*"This movie has good premises. Looks like it has a nice plot, an exceptional cast, first class actors and Stallone gives his best. But it sucks"*

5 POSITIVE TERMS VS 1 NEGATIVE

## 2) Usare tecniche supervisionate no NLP

# Come analizzare i social?



## 3) Stimare la distribuzione aggregata

# Come analizzare i social?



## 4) Non solo "sentiment": analizziamo le opinioni

VOICES
from the Blogs

# Tassonomia delle tecniche di analisi testuale

# Principi della Text Analysis

**Every quantitative linguistic model is wrong, but some can be useful**

# Principi della Text Analysis

**Every quantitative linguistic model is wrong, but some can be useful**

**Quantitative methods help, but cannot replace human**

# Principi della Text Analysis

**Every quantitative linguistic model is wrong, but some can be useful**

**Quantitative methods help, but cannot replace human**

**There exists not BEST or IDEAL technique of text analysis**

# Principi della Text Analysis

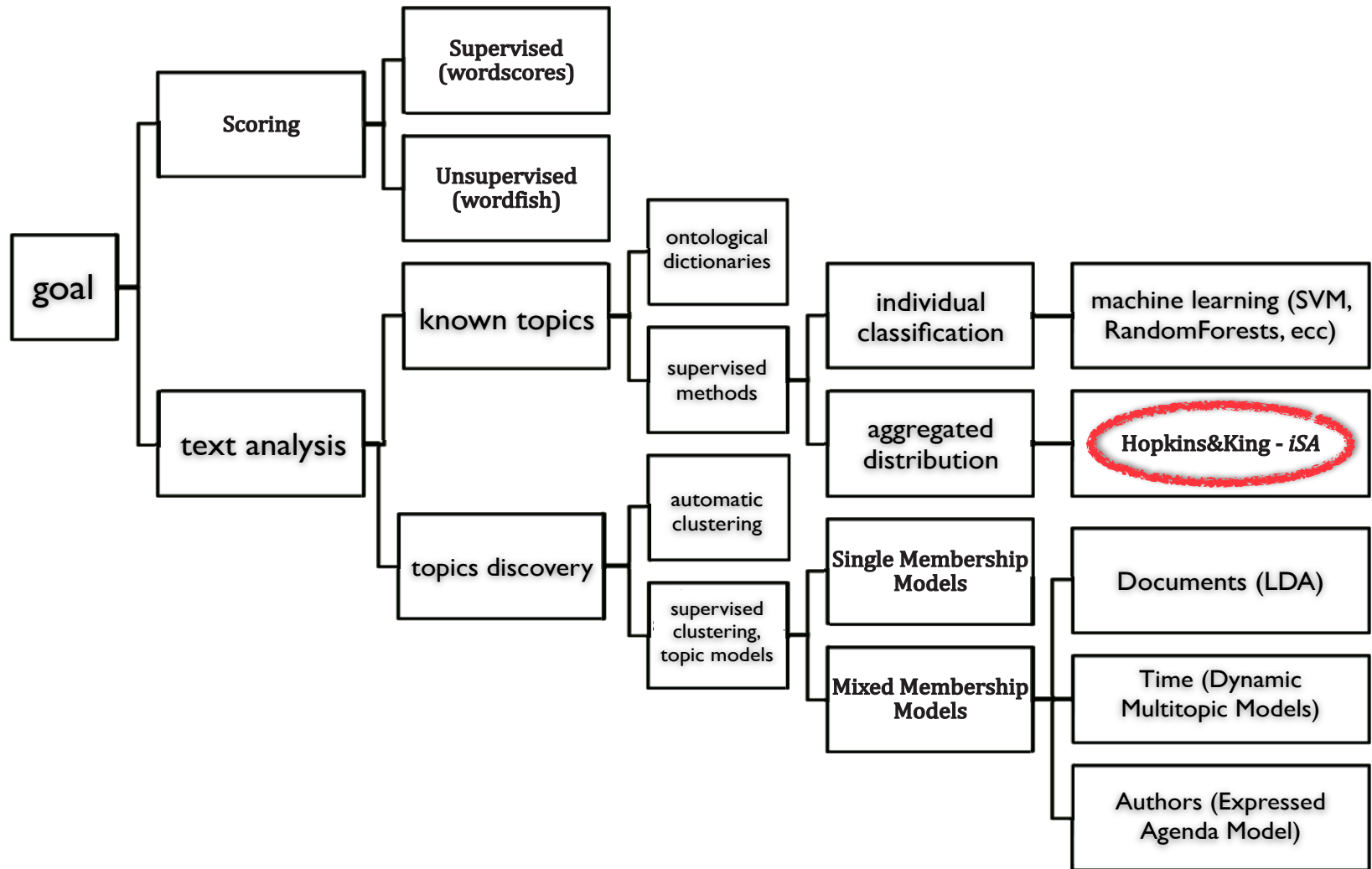**Every quantitative linguistic model is wrong, but some can be useful**

**Quantitative methods help, but cannot replace human**

**There exists not BEST or IDEAL technique of text analysis**

**Validate your analysis**

# Tassonomia delle tecniche

# L'innovazione *i*SA ®

❖ La tecnologia iSA® (*integrated Sentiment Analysis*) sviluppata da **VOICES** rende possibile studiare i Big Data con la profondità di una analisi **qualitativa**

❖ iSA® è il **migliore algoritmo esistente al mondo** per efficacia, velocità di analisi e robustezza nello svolgere analisi sulle opinioni espresse sui Big Data

# L'innovazione *iSA* ®

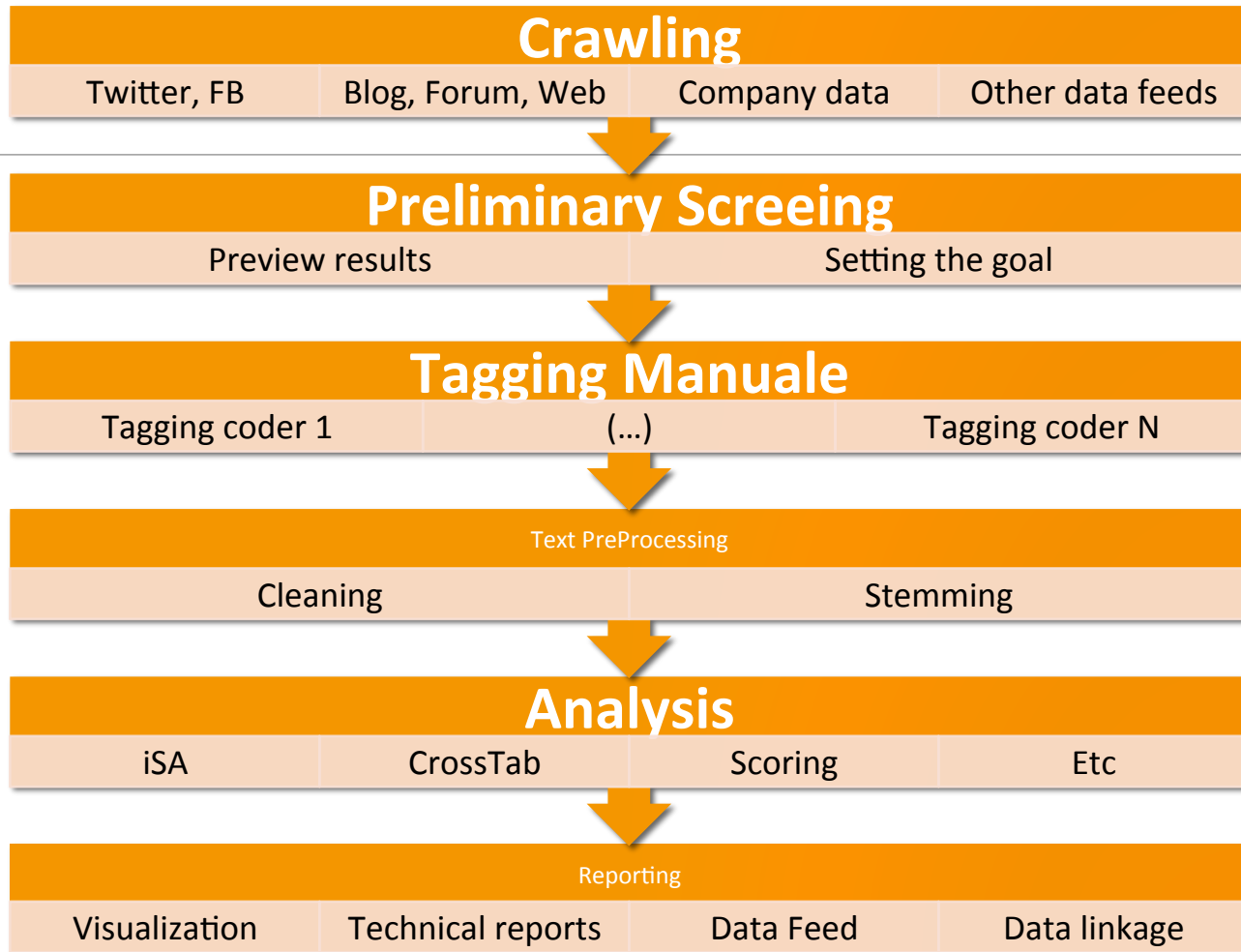| Tool prima generazione | iSA® |
|---|---|
| Volumi di conversazione | **Significato** delle conversazioni |
| Sentiment positivo / negativo | **Motivazioni** dietro le opinioni |
| Conteggio parole | **Analisi** statistica |
| Dipendenti dalla lingua | **Indipendente** dalla lingua |
| Accuratezza medio-bassa (~85%) | **Accuratezza** alta ( > 97%) |

# Come funziona in pratica?

# Come funziona in pratica?

# Il workflow tipico

| Crawling | | | |
|---|---|---|---|
| Twitter, FB | Blog, Forum, Web | Company data | Other data feeds |

| Preliminary Screeing | |
|---|---|
| Preview results | Setting the goal |

| Tagging Manuale | | |
|---|---|---|
| Tagging coder 1 | (…) | Tagging coder N |

| Text PreProcessing | |
|---|---|
| Cleaning | Stemming |

| Analysis | | | |
|---|---|---|---|
| iSA | CrossTab | Scoring | Etc |

| Reporting | | | |
|---|---|---|---|
| Visualization | Technical reports | Data Feed | Data linkage |

# Come funziona in pratica?

post#1: "il nucleare conviene perché è economico"
post#2: "il nucleare produce scorie"
post#3: "il nucleare mi fa paura per le radiazioni, le scorie e non riduce l'inquinamento"

Si dividono i post dei social network in due gruppi:

✳il **train set** : ovvero i testi che verranno letti da codificatori umani 👤 e con i quali si istruisce l'agoritmo 🐻 affinché possa eseguire le stime

✳il **test set**: ovvero i testi che non verranno letti dai codificatori ma con i quali l'algoritmo 🐻 stimerà la distribuzione aggregata delle opinioni

Ogni testo, sia del train che del test set, viene scomposto in stilemi/parole detti "stem" dall'algoritmo 🐻

# IsA at work

Post1: Nuclear energy is convenient as it is cheaper

Post2: Nuclear energy produces waste

Post3: Nuclear scarry me bacause of radiation

Training Set: produced by humans (annotators), used to train the model

Test Set: used to estimate the distribution of the opinion

post#1: "il nucleare conviene perché è economico"
post#2: "il nucleare produce scorie"
post#3: "il nucleare mi fa paura per le radiazioni, le scorie e non riduce l'inquinamento"

Codifica manuale                    Stemming

| | Post | Di | Word: nucleare | Word: paura | Word: radiazioni | Word: inquinamento | Word: scorie | Word: economico |
|---|---|---|---|---|---|---|---|---|
| train set | post#1 | a favore | 1 | 0 | 0 | 0 | 0 | 1 |
| test set | post#2 | NA | 1 | 0 | 0 | 0 | 1 | 0 |
| train set | post#3 | contro | 1 | 1 | 1 | 1 | 1 | 0 |
| train set | post#4 | contro | 1 | 1 | 1 | 1 | 1 | 0 |
| train set | post#5 | a favore | 1 | 0 | 1 | 0 | 0 | 1 |
| | ... | ... | ... | ... | ... | ... | ... | ... |
| test set | post#1000 | NA | 1 | 0 | 0 | 0 | 0 | 1 |

N

| Post | Di | Word: nuclearE | Word: paura | Word: radiazioni | Word: inquinamento | Word: scorie | Word: economico |
|------|-----|-----|-----|-----|-----|-----|-----|
| post#1 | a favore | 1 | 0 | 0 | 0 | 0 | 1 |

*train set*

Di = "a favore"    Si = (1,0,0,0,0,1)

**Goal**: stima della distribuzione **P(D)**

VOICES
from the Blogs

| | Post | Di | Word: nucleare | Word: paura | Word: radiazioni | Word: inquinamento | Word: scorie | Word: economico |
|---|---|---|---|---|---|---|---|---|
| train set | post#1 | a favore | 1 | 0 | 0 | 0 | 0 | 1 |
| test set | post#2 | NA | 1 | 0 | 0 | 0 | 1 | 0 |
| train set | post#3 | contro | 1 | 1 | 1 | 1 | 1 | 0 |
| train set | post#4 | contro | 1 | 1 | 1 | 1 | 1 | 0 |
| train set | post#5 | a favore | 1 | 0 | 1 | 0 | 0 | 1 |
| | ... | ... | ... | ... | ... | ... | ... | ... |
| test set | post#1000 | NA | 1 | 0 | 0 | 0 | 0 | 1 |

Dizionario Treccani (italiano): 270k lemmi
Oxford Dictionary (English): 650k lemmi

In realtà, per ciascun argomento nel linguaggio comune si tende ad utilizzare al massimo M = 200 o 500 "stilemi" e questo rende possibile effettuare l'analisi statistica

L'analisi resta difficile da approcciare perché avremo potenzialmente $2^M$ righe diverse composte da 0 e 1 (ovvero tra $1,6*10^{50}$ e $3,3*10^{150}$)

VOICES
from the Blogs

Approcc **goal** **train & test** **train+test**

$$P(D) = P(D|S) * P(S)$$

modello statistico classico
produce missclassification

distribuzione degli stem

Stime distorte
alta variabilità

**All machine learning methods affected**, choose your own and add to the list below:
- Support Vector Machines
- Random Forests
- Neural Networks
- ecc

Problemi di tipo statistico: questo approccio cerca di rintracciare le opinioni considerando tutti i 2^M possibili "stilemi" → time-consuming e non garantisce risultati accurati perché diventa improbabile rintracciare "stilemi" che stiano davvero esprimendo un'opinione

Met**Approccio statistico innovativo** (King&Hopkins, 2010)

*train+test*   *train*   *goal*

$$P(S) = P(S|D) * P(D)$$

$$P(S|D)^{-1} * P(S) = P(D)$$

Si guarda alla distribuzione degli Stem in ciascuna categoria di "opinioni" e non il contrario

"Semplice" quanto invertire una matrice!

Nessun problema a gestire la quantità di "Big" Data

VOICES
from the Blogs

Met**Approccio statistico innovativo** (King&Hopkins, 2010)

*train+test*      *train*      *goal*

$$P(S) = P(S|D) * P(D)$$

$$P(S|D)^{-1} * P(S) = P(D)$$

Si guarda alla distribuzione degli Stem in ciascuna categoria di "opinioni" e non il contrario

"Semplice" quanto invertire una matrice!

Nessun problema a gestire la quantità di "Big" Data

The space "Opinion x Stems" = $\mathcal{D} \times \bar{\mathcal{S}}$

# Tipica performance statistica di iSA vs competitors

## "Bias"

# Tipica performance statistica di iSA vs competitors

## "Precision" (Large Movie dataset)

# Forecasting Peruvian Elections with twitter

# Sentiment-enhanced Multidimensional Analysis of Online Social Networks:

Perception of the Mediterranean Refugees Crisis

Mauro Coletto ∗†, Andrea Esuli †, Claudio Lucchese †, **Cristina Ioana Muntean** †, Franco Maria Nardini †, Raffaele Perego †, Chiara Renso †

∗ IMT School for Advanced Studies Lucca - ITALY
† ISTI - CNR Pisa - ITALY

# Refugees Crisis Perception Analysis

AQ1: What is the evolution of the discussions about refugees migration in Twitter?

AQ2: What is the sentiment of users across Europe in relation to the refugee crisis? What is the evolution of the perception in countries affected by the phenomenon?

AQ3: Are users more polarized in countries most impacted by the migration flow?

# Analytical Framework

An analytical framework to interpret social trends from large tweet collections by extracting and crossing information about the following three dimensions:

◦ Time

◦ Space
  ◦ User location
  ◦ Location mentions

◦ Sentiment
  ◦ Tweet sentiment
  ◦ User sentiment

| Symbol | Description | # Total |
|--------|-------------|---------|
| $\mathcal{G}$ | Collected English tweets | 97,693,321 |
| $\mathcal{T}$ | Tweets related to the refugee crisis | 1,238,921 |
| $\mathcal{T}_{c+}$ | Positive sentiment tweets | 459,544 |
| $\mathcal{T}_{c-}$ | Negative sentiment tweets | 387,374 |
| $\mathcal{T}_{ML}$ | Tweets with mentioned location | 421,512 |
| $\mathcal{T}_{UL}$ | Tweets with user location | 101,765 |
| $\mathcal{U}$ | Users | 480,660 |
| $\mathcal{U}_{c+}$ | Users with positive sentiment | 213,920 |
| $\mathcal{U}_{c-}$ | Users with negative sentiment | 104,126 |
| $\mathcal{U}_{L}$ | Users with country location | 47,824 |

Perform multidimensional analyses considering content and locations in time

# Deriving Sentiment

HASHTAG CLASSIFICATION

Initial seeds

#refugeeswelcome

#refugessnotwelcome

Positive Hashtags
Negative Hashtags

Enrich hashtag seeds
from #-tag co-occurrence

USER CLASSIFICATION

Positive Users
Negative Users

Positive Tweets
Negative Tweets

TWEET CLASSIFICATION

# Sentiment on migration topics: Perception of the Mediterranean Refugee Crisis

**European country mentions**

**Africa & Middle East country mentions**

# Sentiment on migration topics: Perception of the Mediterranean Refugee Crisis

- **Internal and external perception by country**
  - Index **ρ** - the ratio between pro refugees users and against refugees users
  - Red means a higher predominance of positive sentiment, higher ρ
  - Yellow means a higher predominance of negative sentiment, lower ρ



(a) Global perception     (b) Internal perception     (c) External perception

# Sentiment Analysis in UK

Positive and negative users for different cities in UK before and after September 4 (death of Alan Kurdi, borders between AT HU, Germany welcomes refugees).

◦ bars show the number of polarized positive and negative users by city
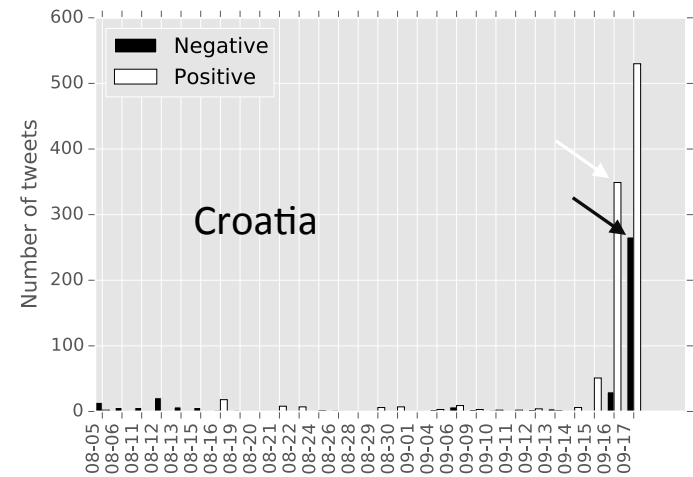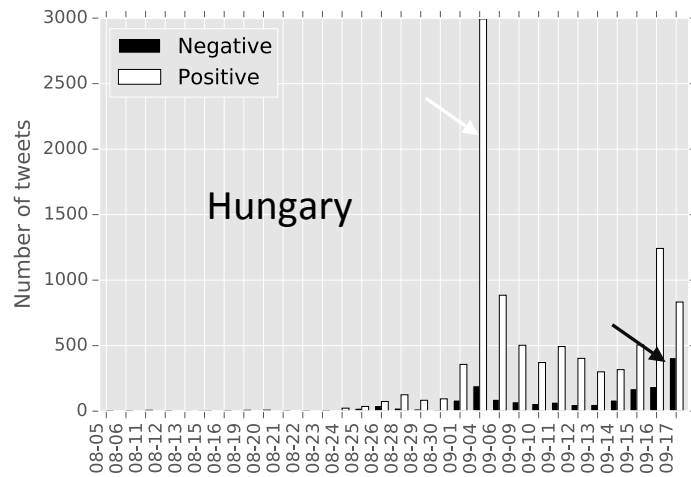
◦ the heat map in background indicates the value of ρ

# Sentiment Analysis
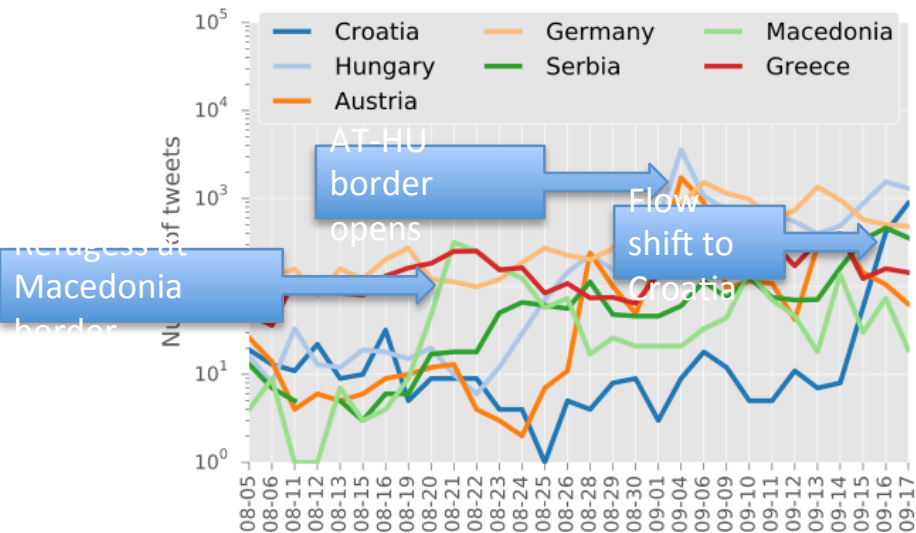For mentioned locations analysis and tweet sentiment

# Sentiment Analysis
For mentioned locations analysis and tweet sentiment

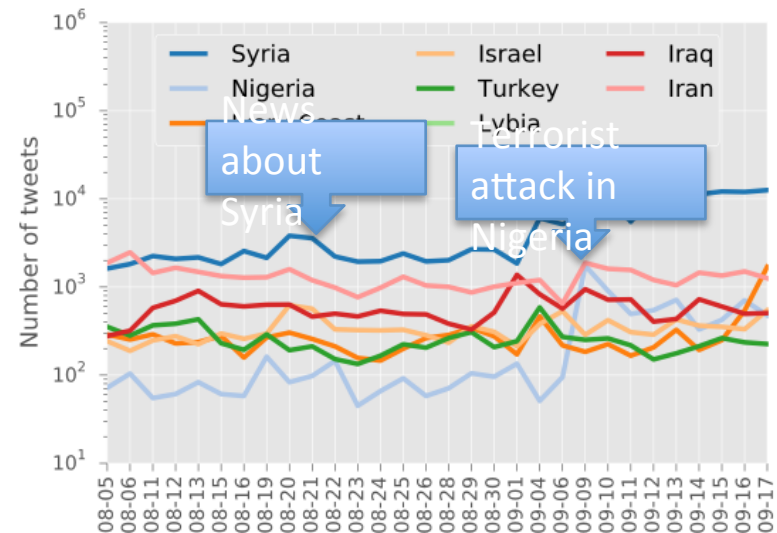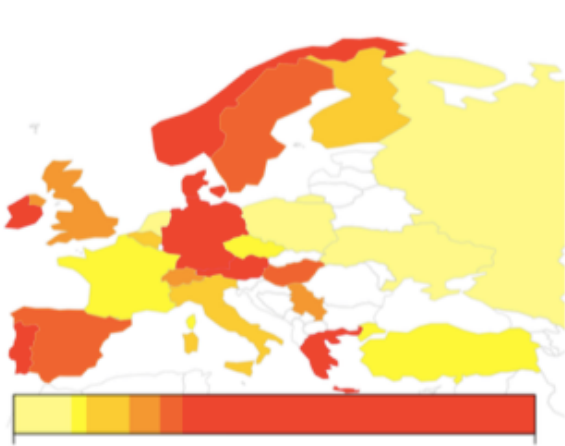# Space and Time analysis



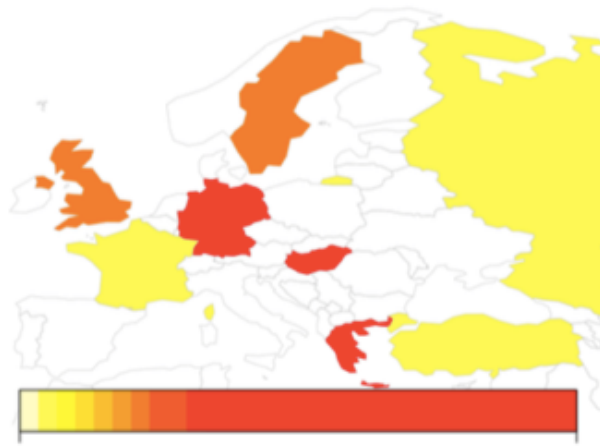**European country mentions**  **Africa & Middle East country mentions**

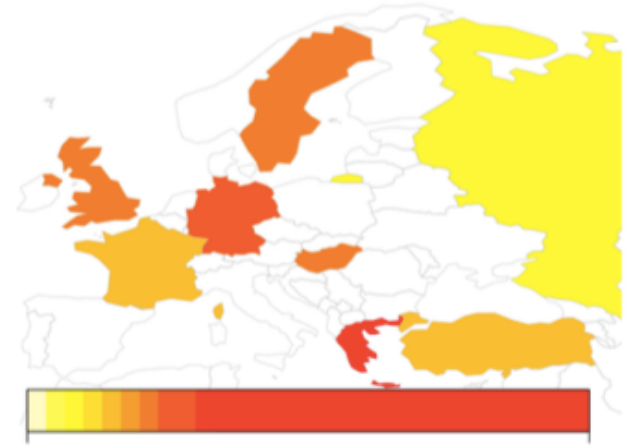# Sentiment Analysis

Internal and external perception by country

◦ Index **ρ** - the ratio between pro refugees users and against refugees users

◦ Red means a higher predominance of positive sentiment, higher ρ



(a) Overall.



(b) Internal perception.



(c) External perception.

THANK YOU !

Questions?