

# Social media **fingerprints** of *unemployment*

---

Alejandro Llorente (IIC, Spain)

Manuel García Herranz (UAM and United Nations, NY)

Manuel Cebrián (NICTA, Australia)

Esteban Moro (UC3M, Spain)

WE  
*are*  
WHAT WE  
*repeatedly*  
DO.

-aristotle

# Computational Social Science

**You are** what you repeatedly do [Aristóteles]

Using BigData to infer behavior or society situation

## Situation

Demographics  
Health  
Economy  
Unemployment  
Transportation  
Geography  
Politics

## Behavior

Social  
Mobility  
Activity  
Content

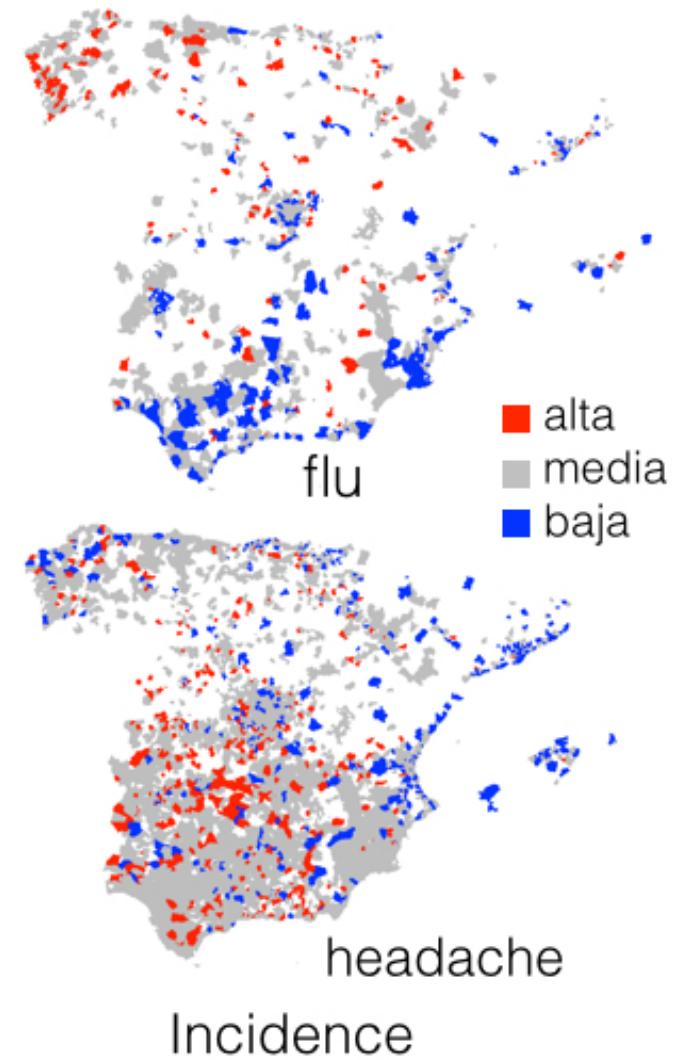
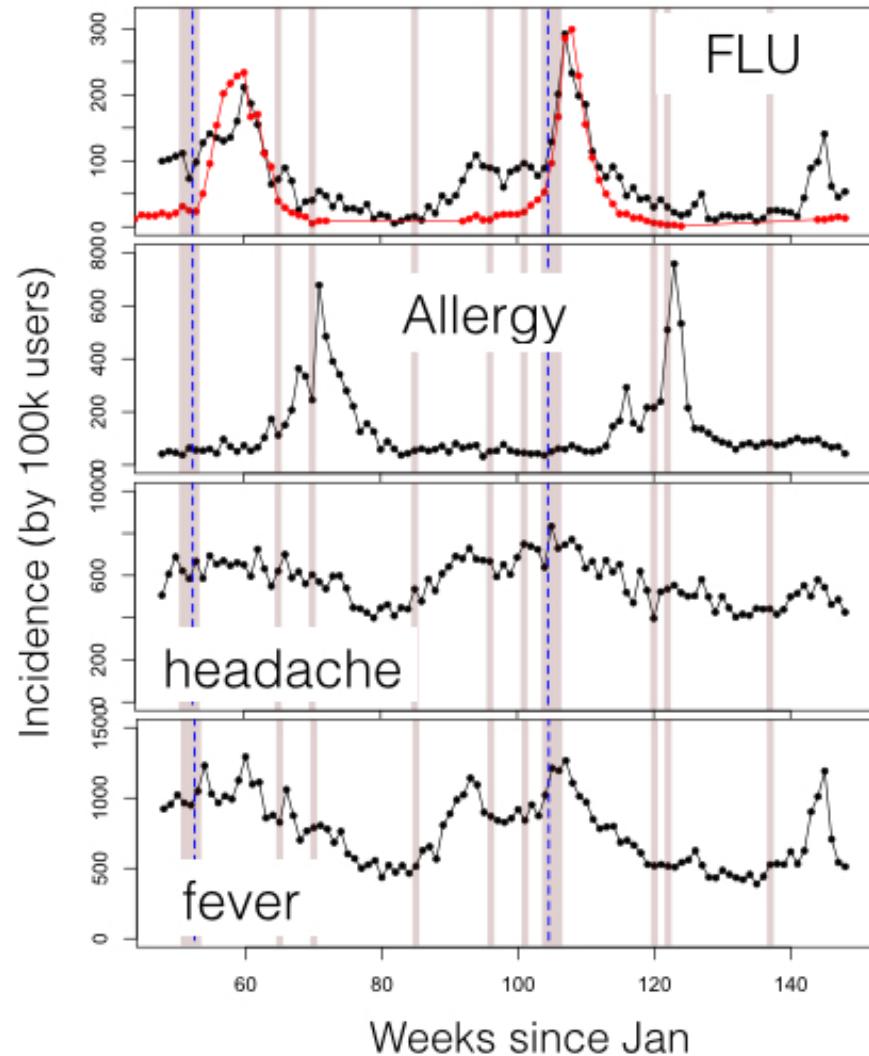
## Observation

Surveys  
Credit card  
Mobile phone  
Social media  
Searches  
...

Individual - Group - City

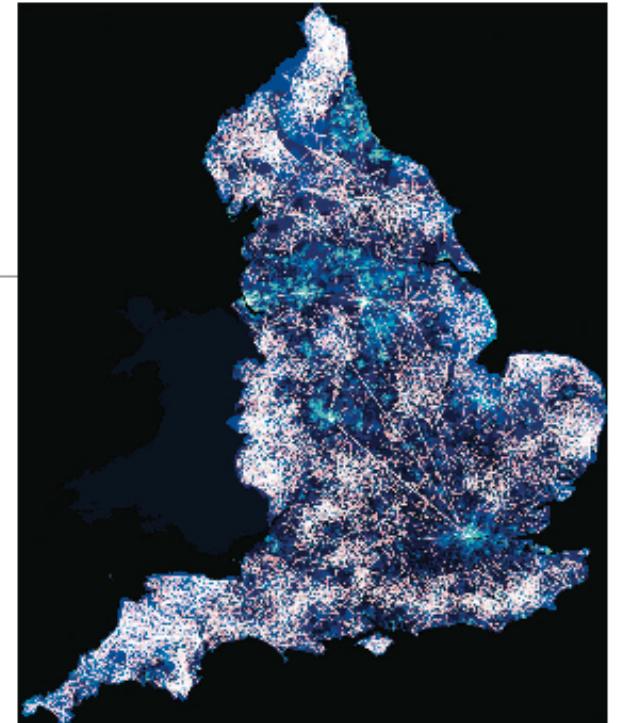
# Computational Social Science

Health -> Content -> Social Media



# Computational Social Science

---



**Economy -> Behavior -> Observation**

Eagle et al, Science 2010

Development -> Social/Mobility (city) diversity -> Mobile phones

Soto, V. et al., Prediction of Socioeconomic Levels using Cell Phone Records 2011.

Socio economical level -> Social/Mobility/Activity -> Mobile phones

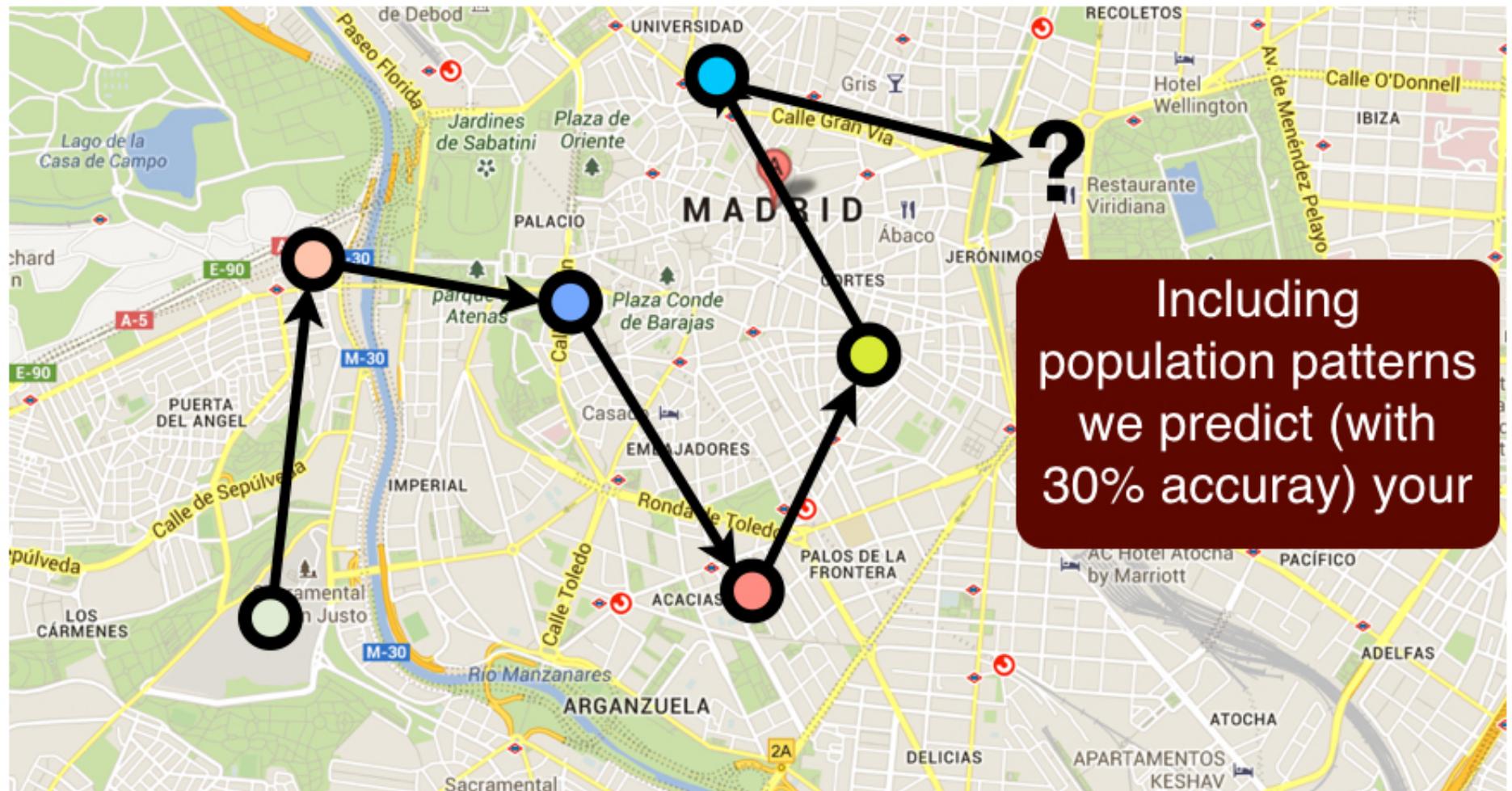
Antenucci, D. et al., 2014. Using Social Media to Measure Labor Market Flows.

Unemployment -> Content -> Social Media

# Computational Social Science

Wealth -> Behavior -> Credit card

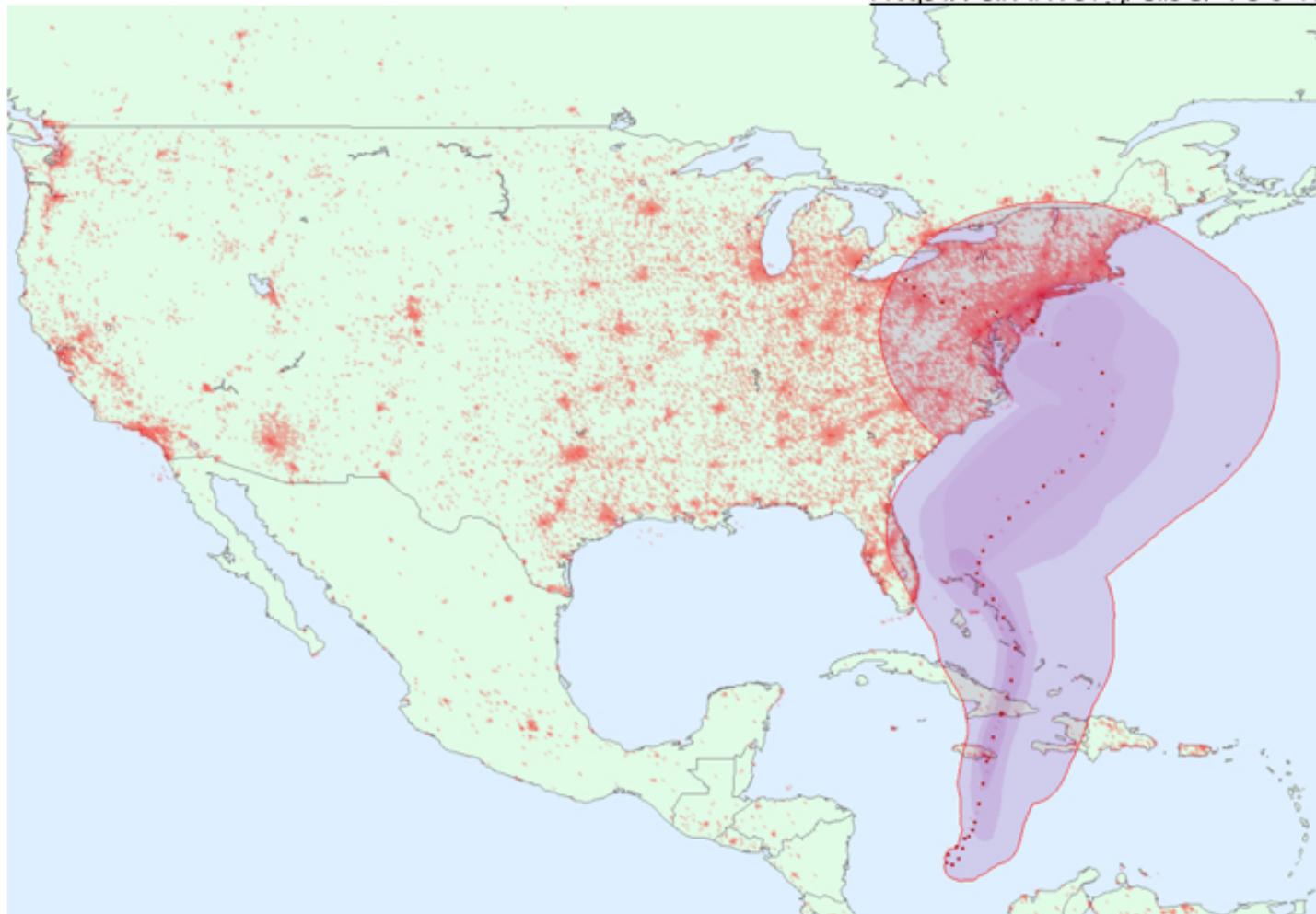
Krumme, C. et al., 2013. The predictability of consumer visitation patterns. Scientific Reports



# Computational Social Science

Disaster -> Behavior -> Social Media

Dataset: 52.55 Million messages, 14 Million users  
Yury Kryvasheyeu, Manuel Cebrián, EM, et al 2015  
<http://arxiv.org/abs/1504.06827>

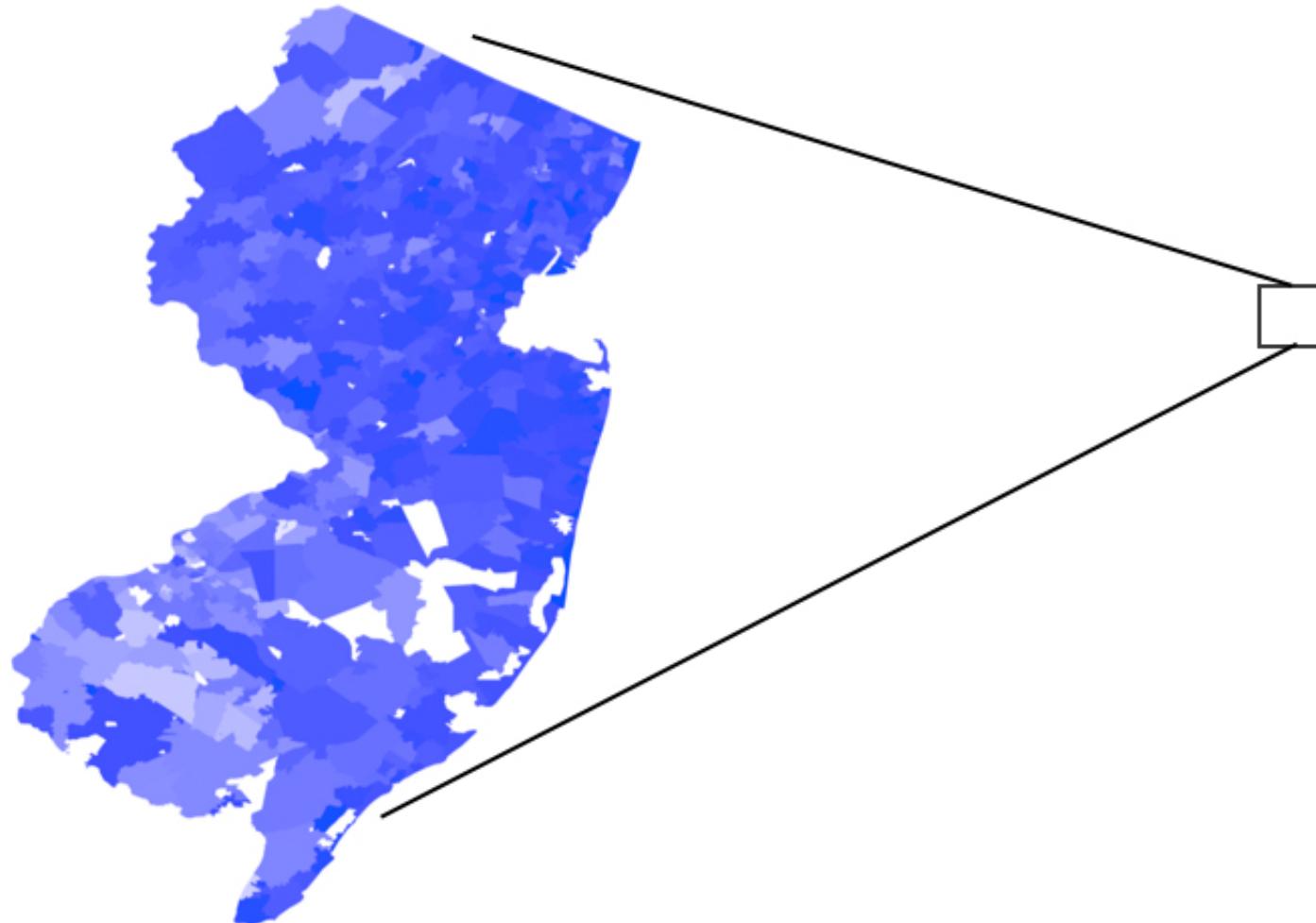


Economical Impact ~ \$1 Billion

# Computational Social Science

Disaster -> Behavior -> Social Media

Dataset: 52.55 Million messages, 14Million users  
Yury Kryvasheyev, Manuel Cebrián, EM, et al 2015  
<http://arxiv.org/abs/1504.06827>

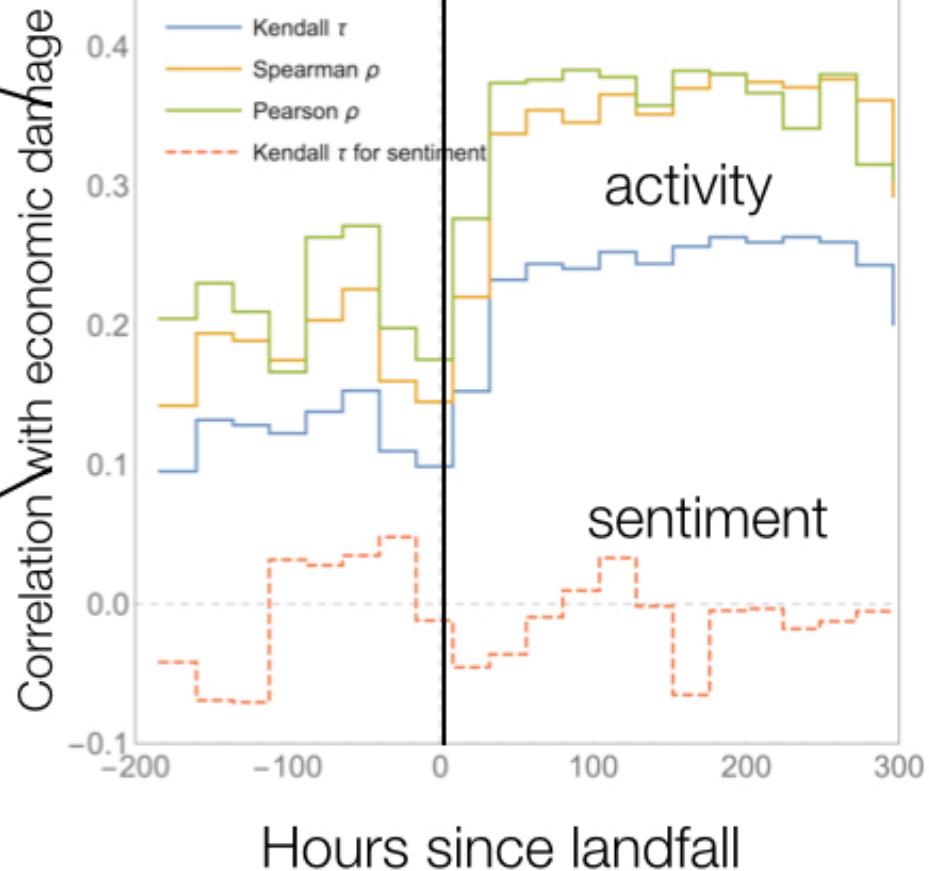
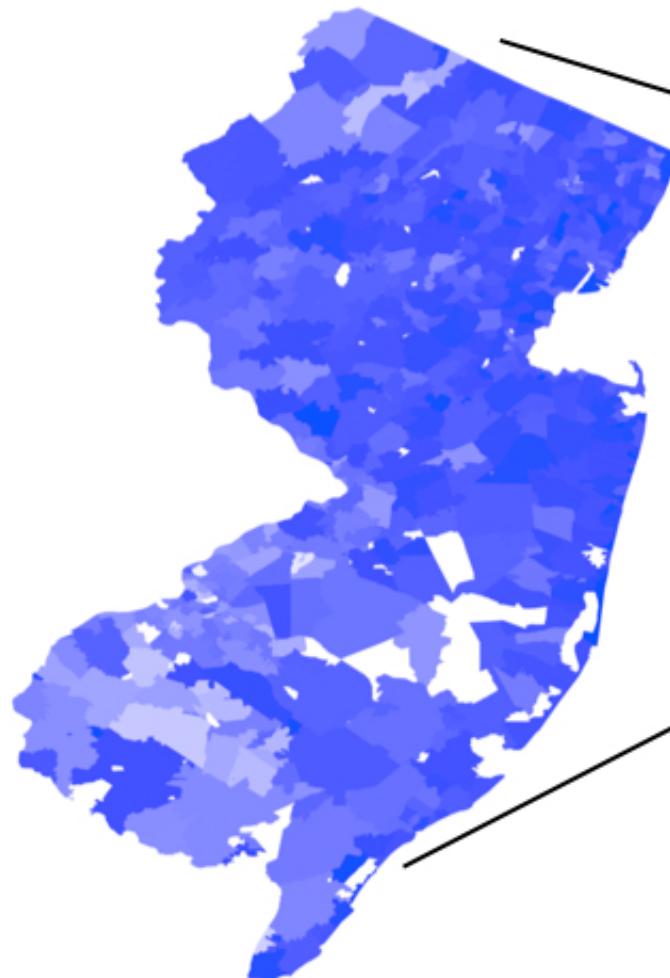


Economical Impact ~ \$1 Billion

# Computational Social Science

Disaster -> Behavior -> Social Media

Dataset: 52.55 Million messages, 14 Million users  
Yury Kryvasheyev, Manuel Cebrián, EM, et al 2015  
<http://arxiv.org/abs/1504.06827>

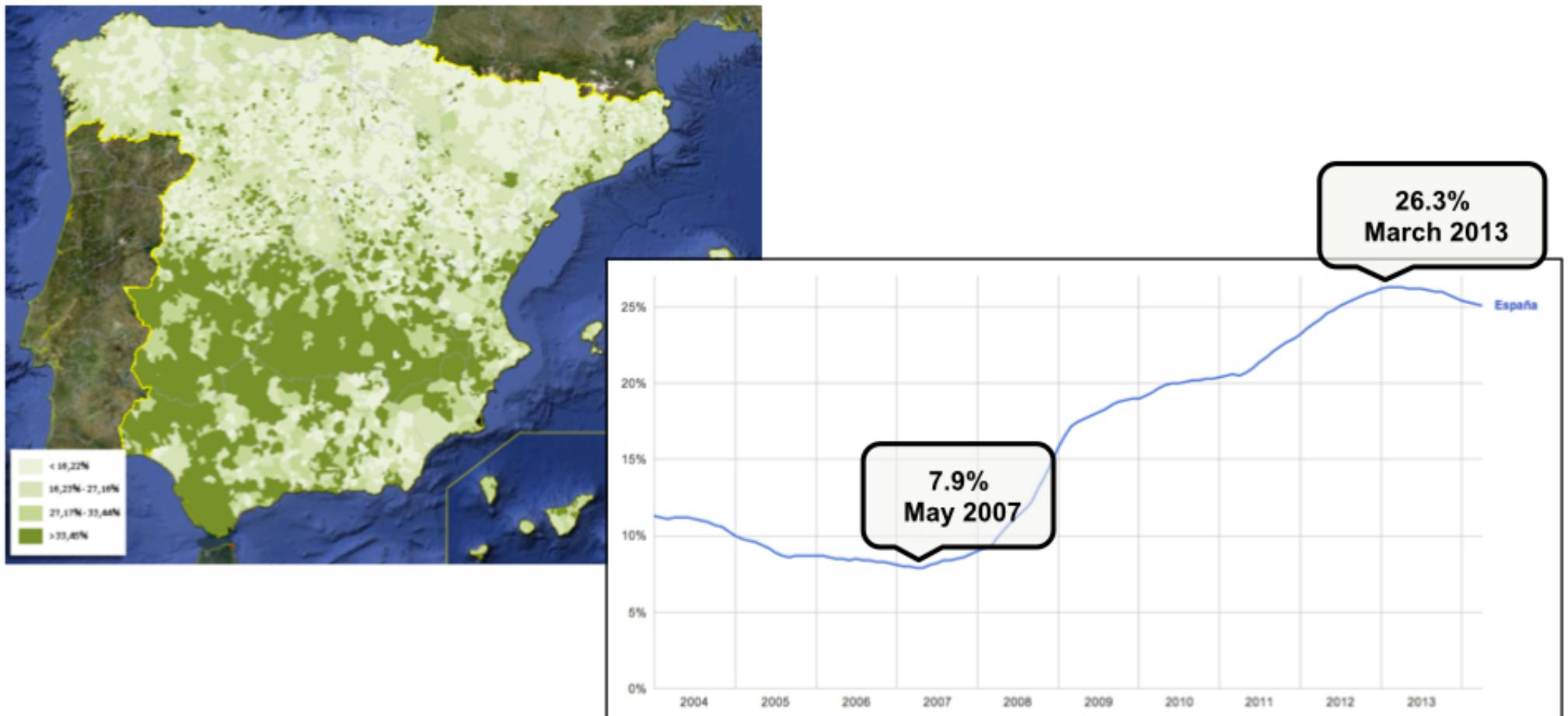


Economical Impact ~ \$1 Billion

# Our objective

**Unemployment -> Behavior -> Social media (Twitter)**

- It's not another economical problem in Spain. It is the BIG problem



# Our objective

---

## **Unemployment -> Behavior -> Social media (Twitter)**

- Twitter allows us to study social/mobility/content behaviors together
  - *But, which behavior(s) are the more relevant with respect to unemployment? (360° view)*
- Twitter is **BigSparseData**: some demographic and geographic areas are not well represented, thus
  - *How much of the observed unemployment can be described by Twitter-detected behaviors?*

FEB  
2014

# SPAIN



**47,370,542**

TOTAL POPULATION



77%

23%

URBAN

RURAL

**33,870,948**

INTERNET USERS



72%

INTERNET PENETRATION

**19,600,000**

ACTIVE FACEBOOK USERS



41%

FACEBOOK PENETRATION

**55,740,000**

ACTIVE MOBILE SUBSCRIPTIONS

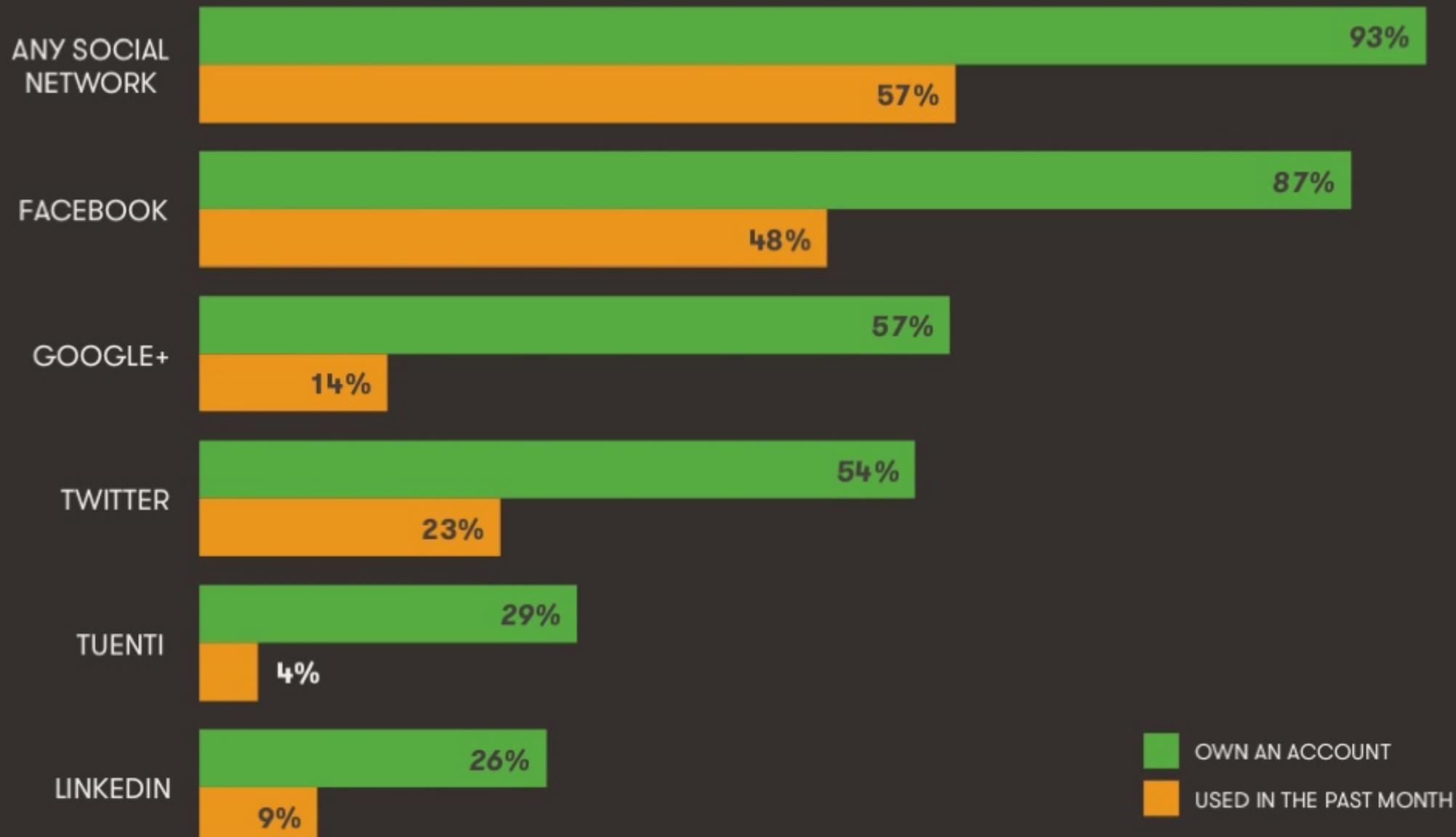


118%

MOBILE SUBSCRIPTION PENETRATION

FEB  
2014

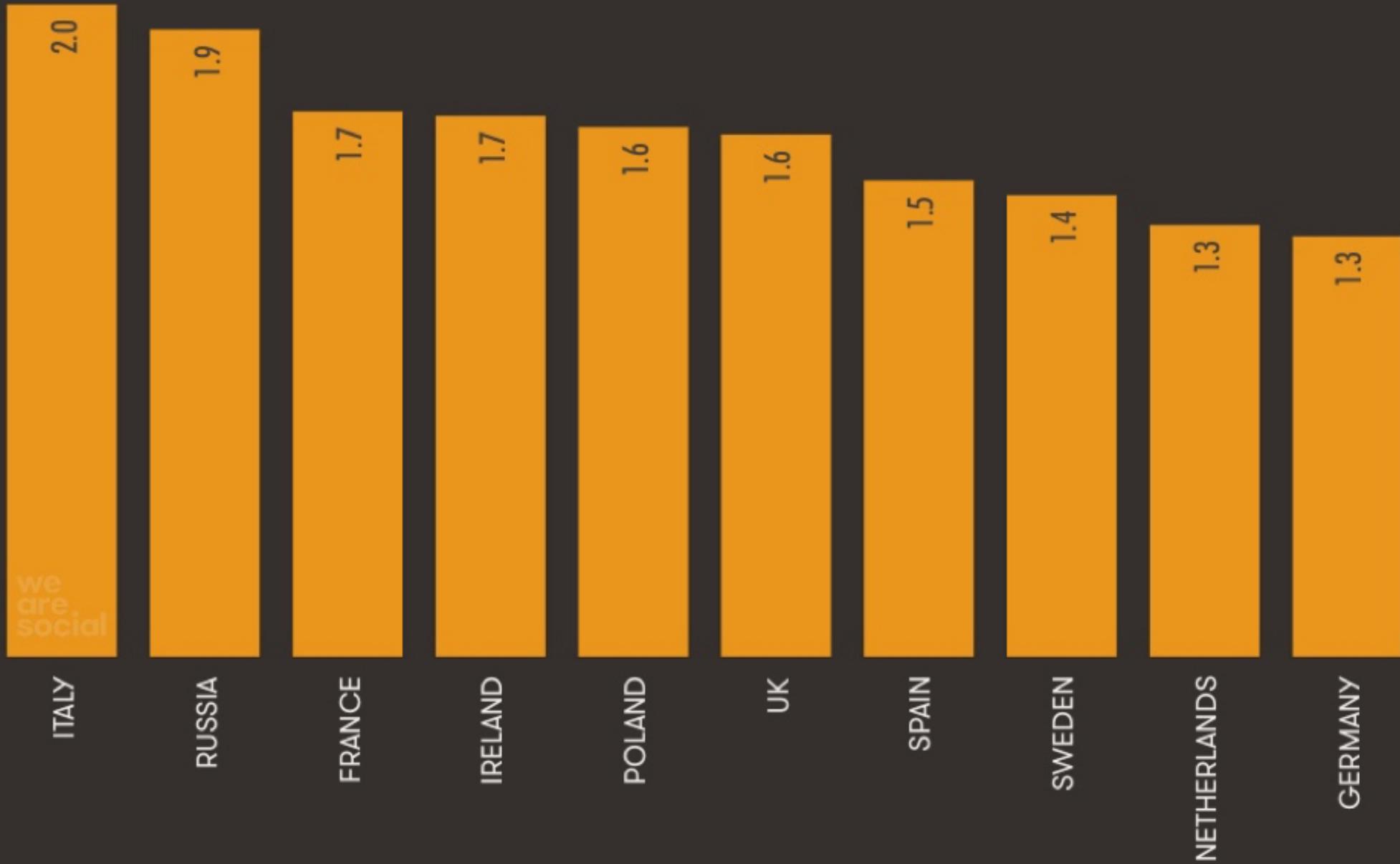
# SPAIN: SOCIAL MEDIA USE



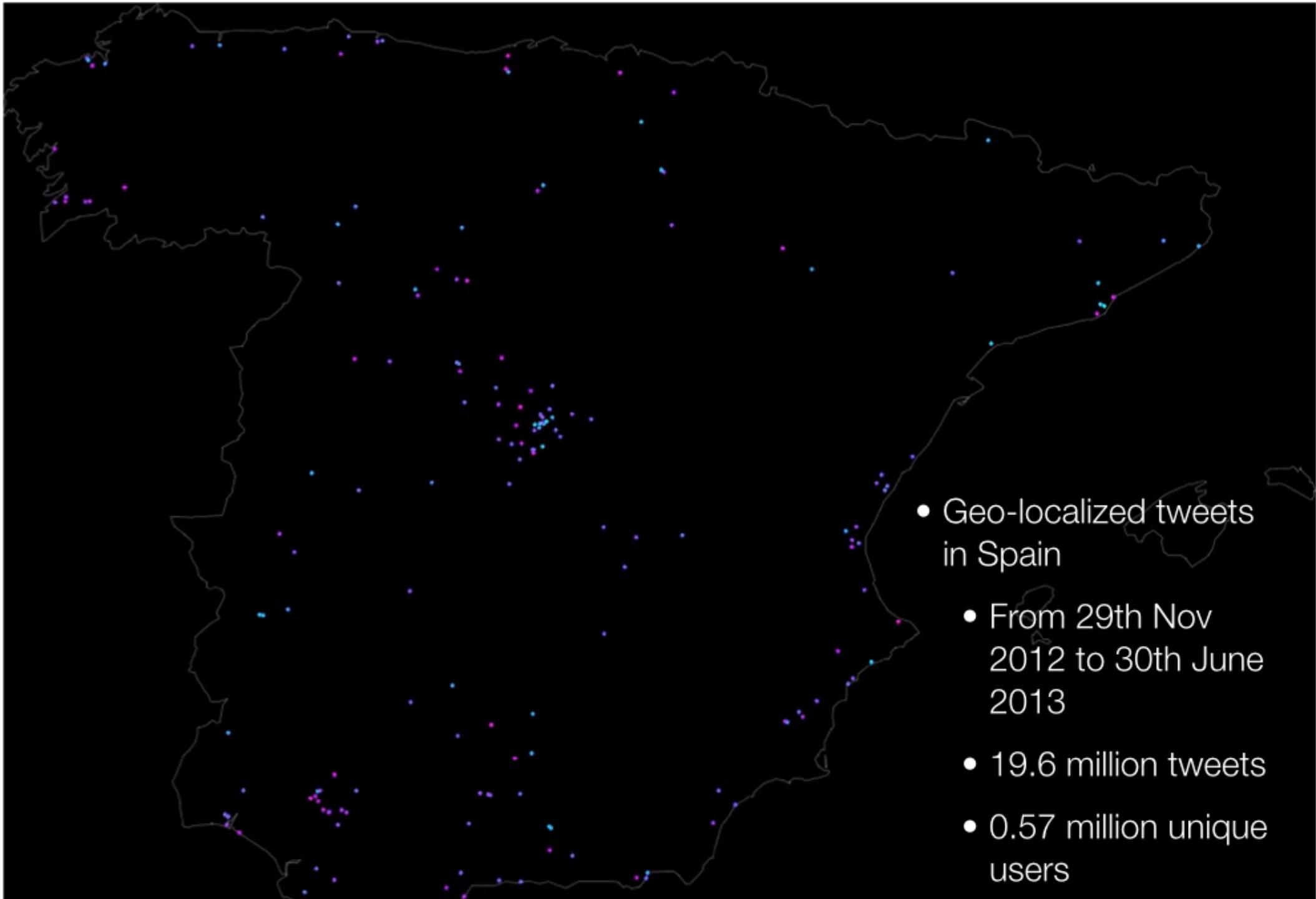
FEB  
2014

# TIME SPENT ON SOCIAL MEDIA

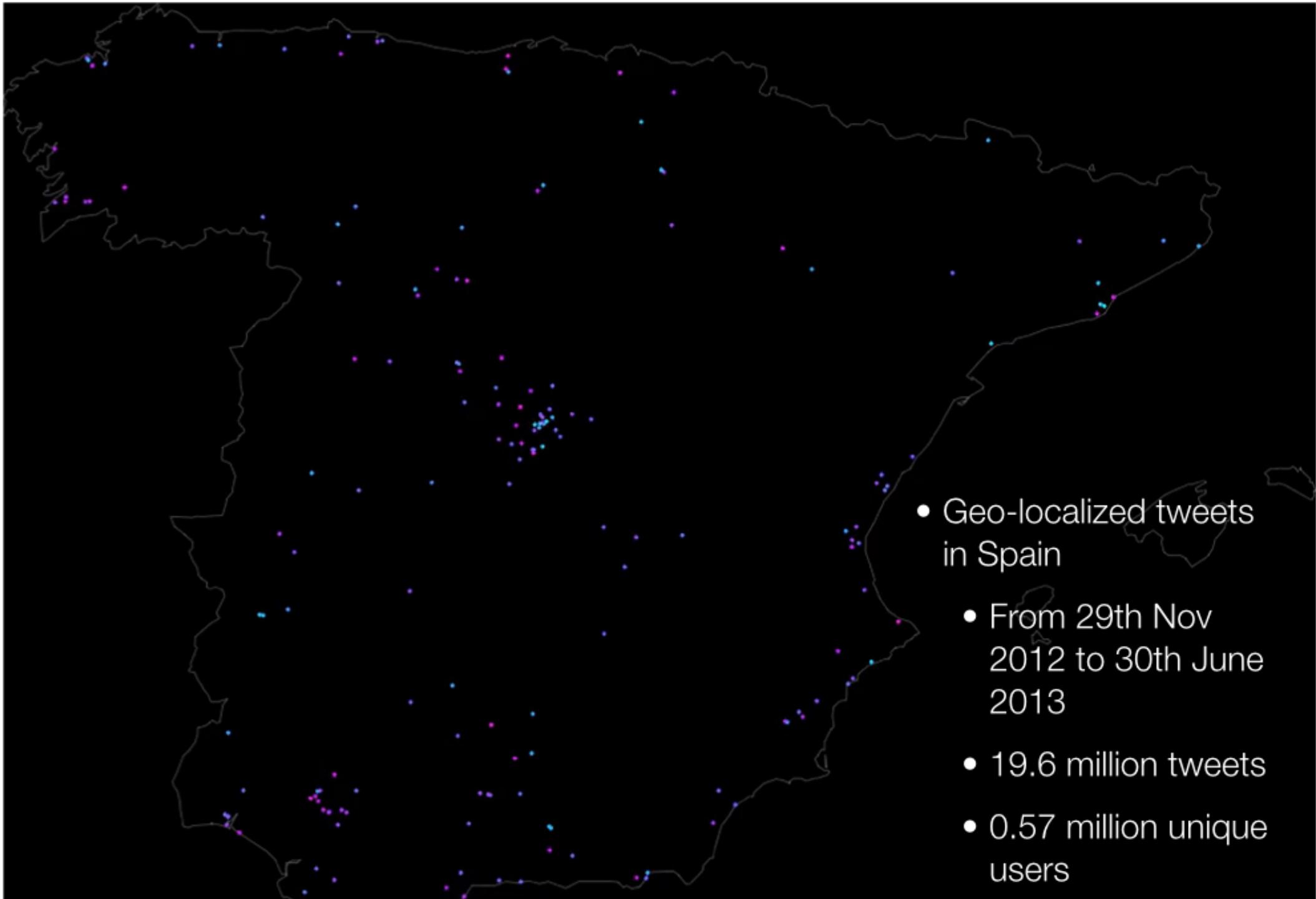
AVERAGE NUMBER OF HOURS PER DAY SPENT BY SOCIAL MEDIA USERS ON ALL SOCIAL CHANNELS



# Our database



# Our database



# Geographical areas in Spain

---

- Municipalities:
  - ~8200 municipalities in Spain
  - Very heterogeneous:  
population ranging from 7 to 3.2
- We propose a functional approach to the definition of the areas based on daily mobility
- *Users' municipality home is where they tweet most*



# Geographical areas in Spain

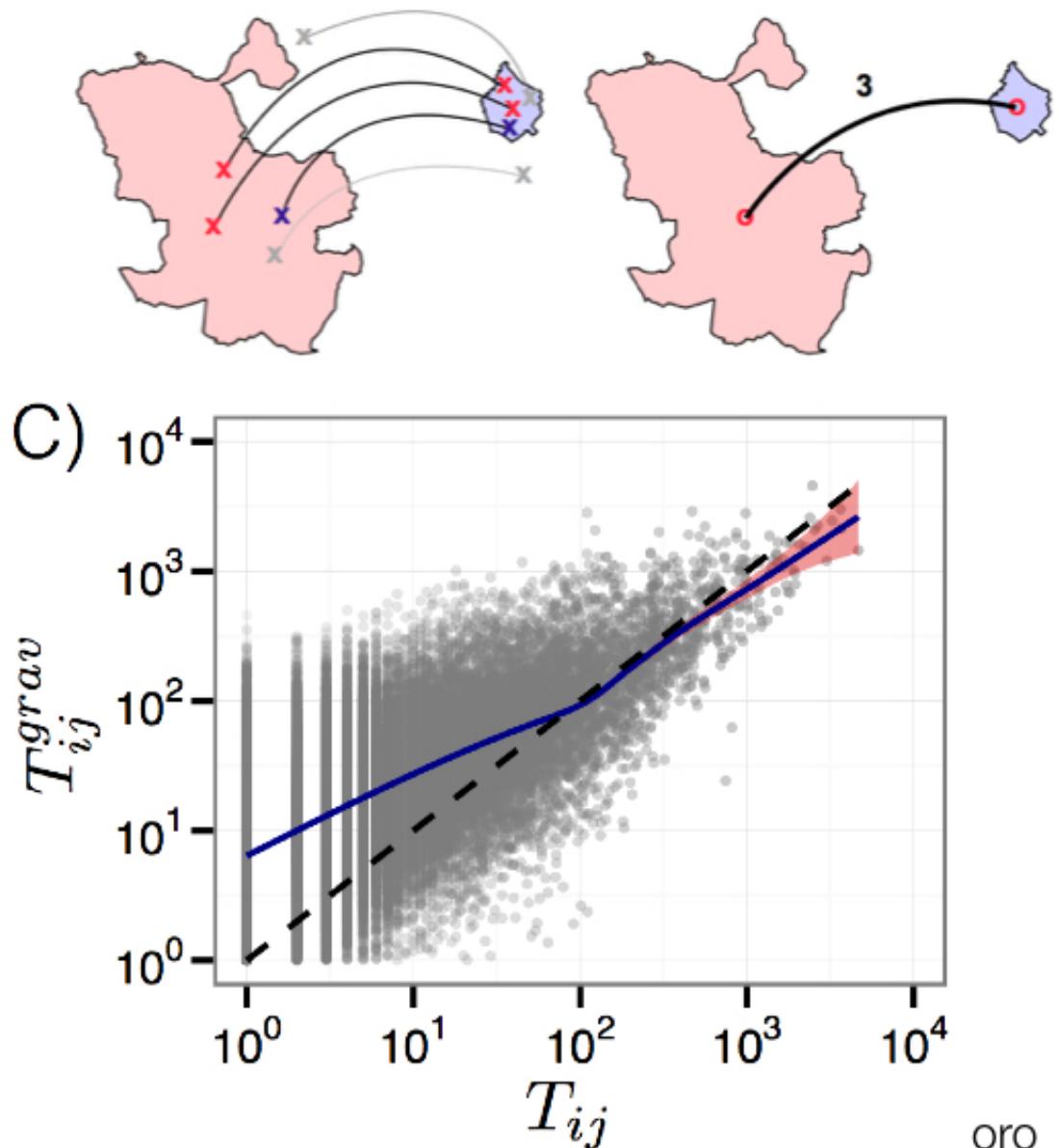
- We take all trips between two municipalities to construct the flow  $T_{ij}$  the number of trips between municipalities
- Flow is well described by the gravity model

$$T_{ij} \simeq T_{ij}^{grav} = \frac{P_i^{\alpha_i} P_j^{\alpha_j}}{d_{ij}^{\beta}}$$

$$\alpha_i \approx \alpha_j = 0.42, \beta = 0.89$$

$$R^2 = 0.69$$

Lenormand, M. et al., 2014. Cross-checking different sources of mobility information.



# (Functional) geographical areas in Spain

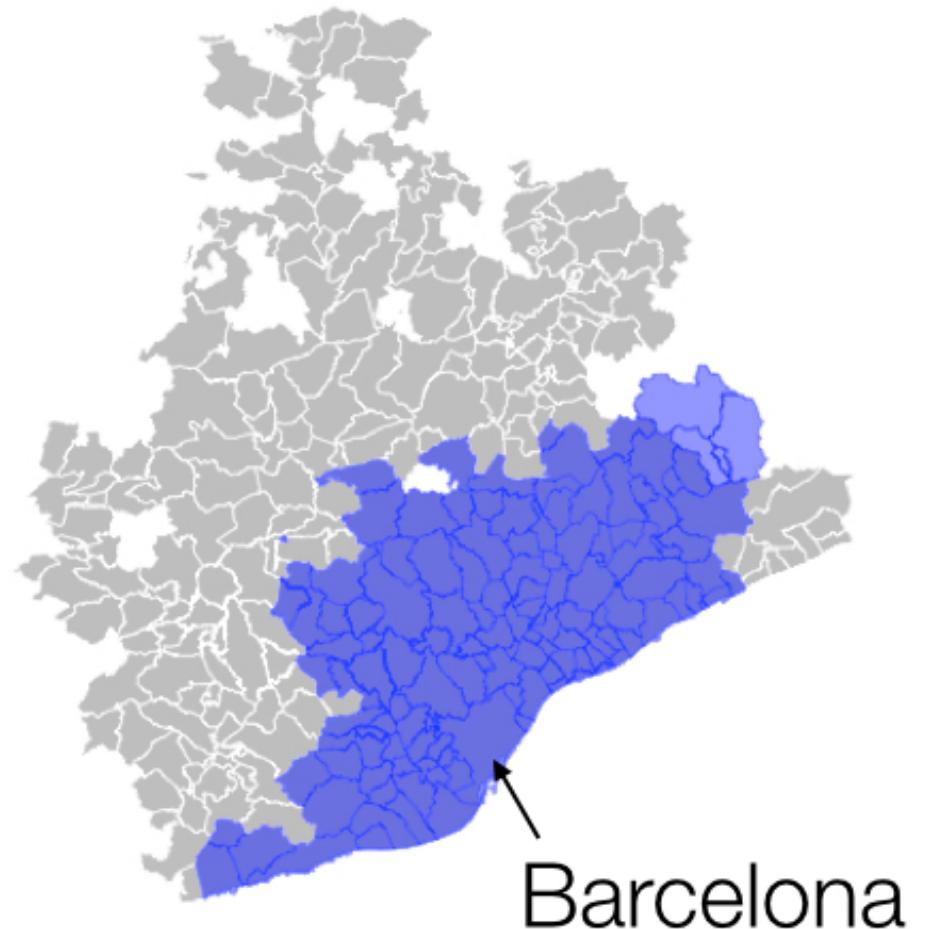
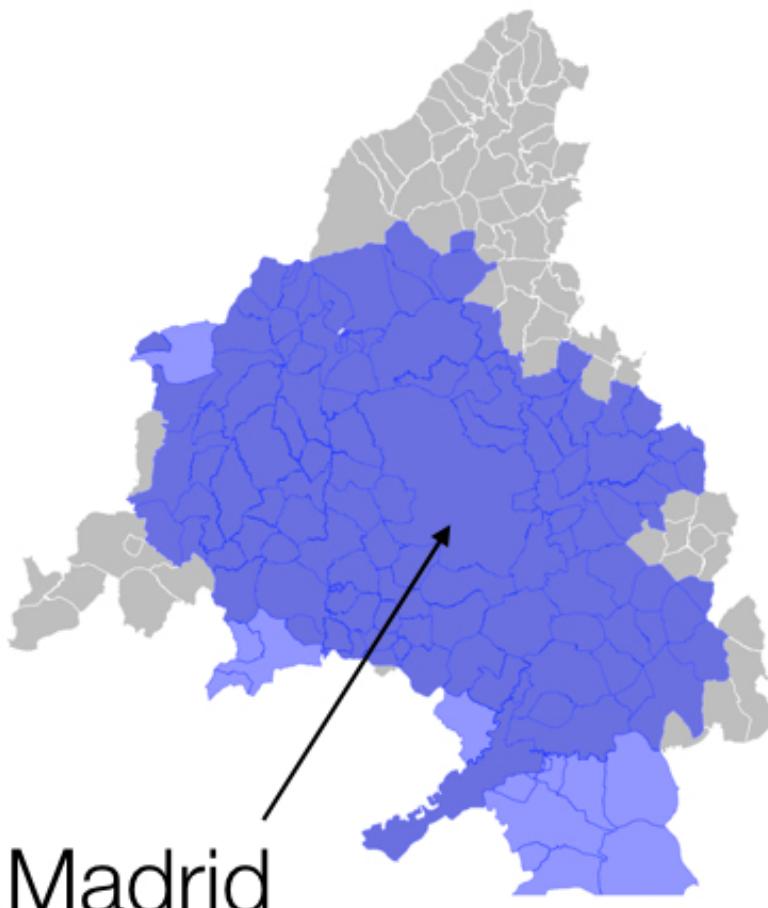
---

- We look for communities in the flow graph
- Infomap algorithm
- 340 functional areas were detected:
  - They are cohesive
  - Statistically robust
  - Modularity is high
  - High overlap with “comarcas”

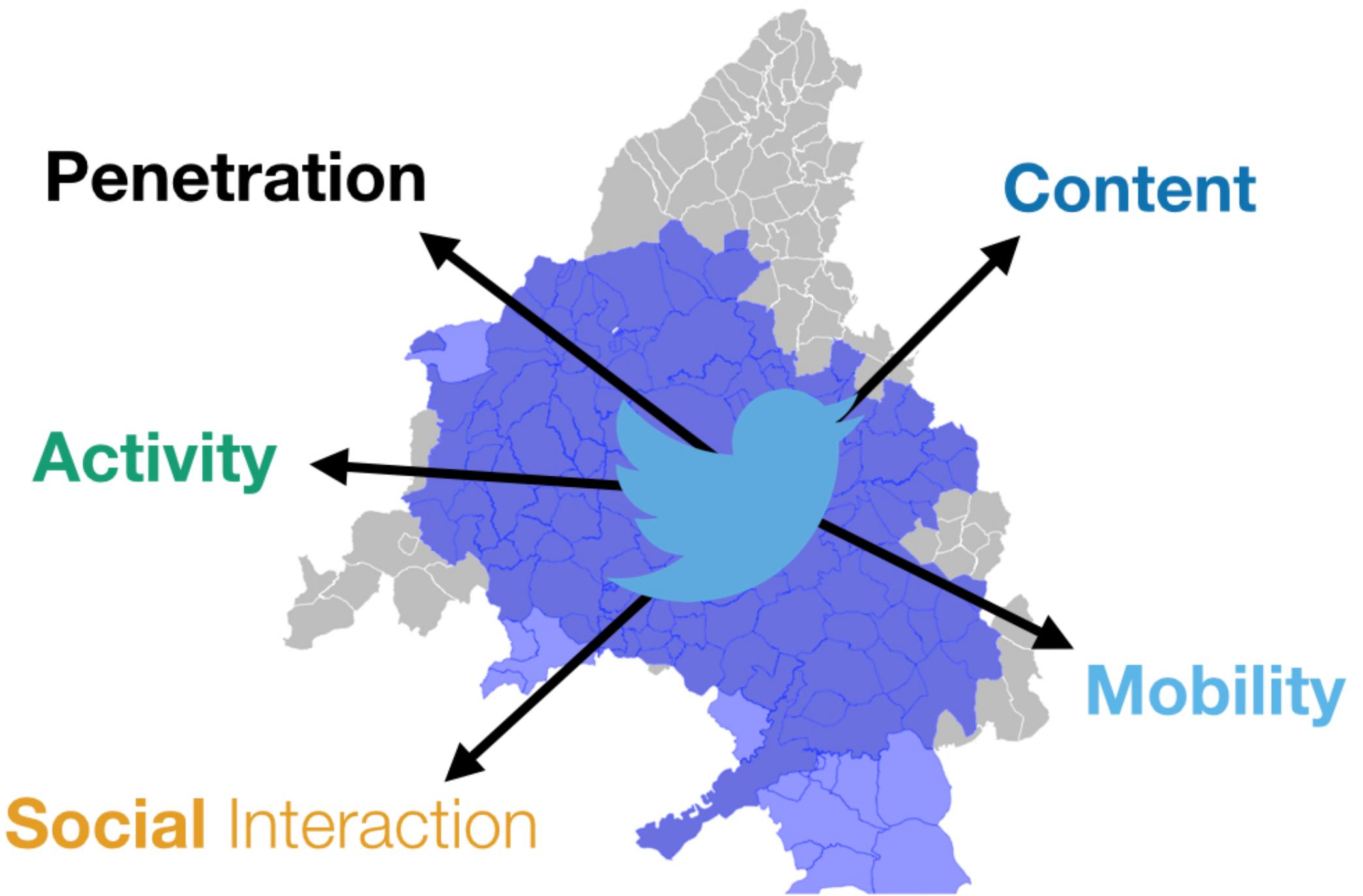


# (Functional) geographic areas in Spain

---



"The piece is absolutely useless, even ridiculous, outside Spain, because the audience cannot hope to understand its significance, nor the performers to play it as it should be played."



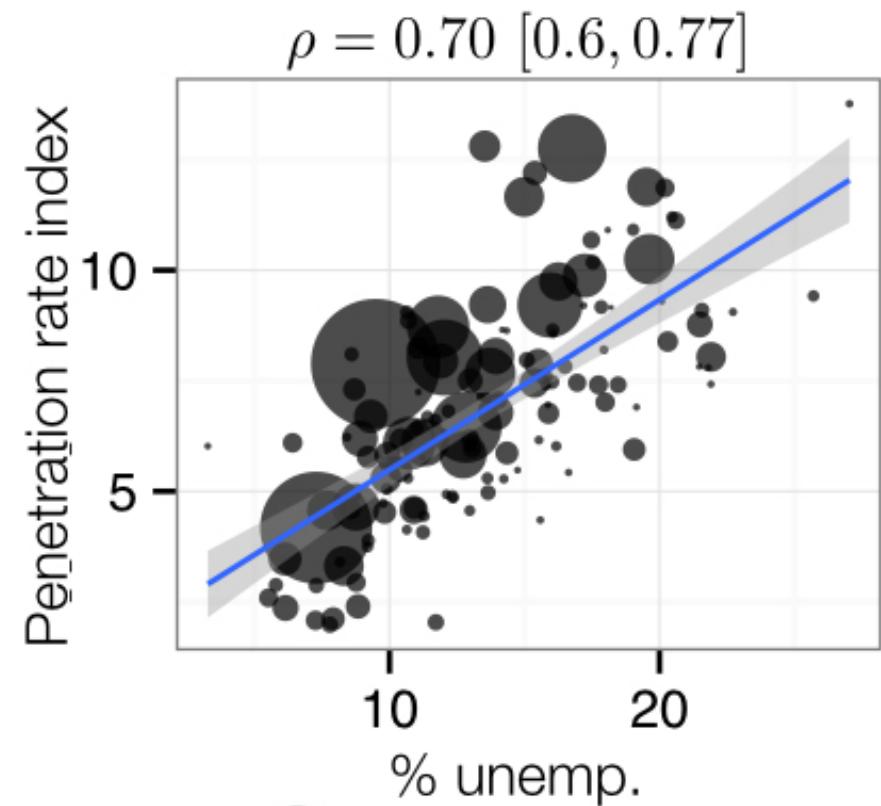
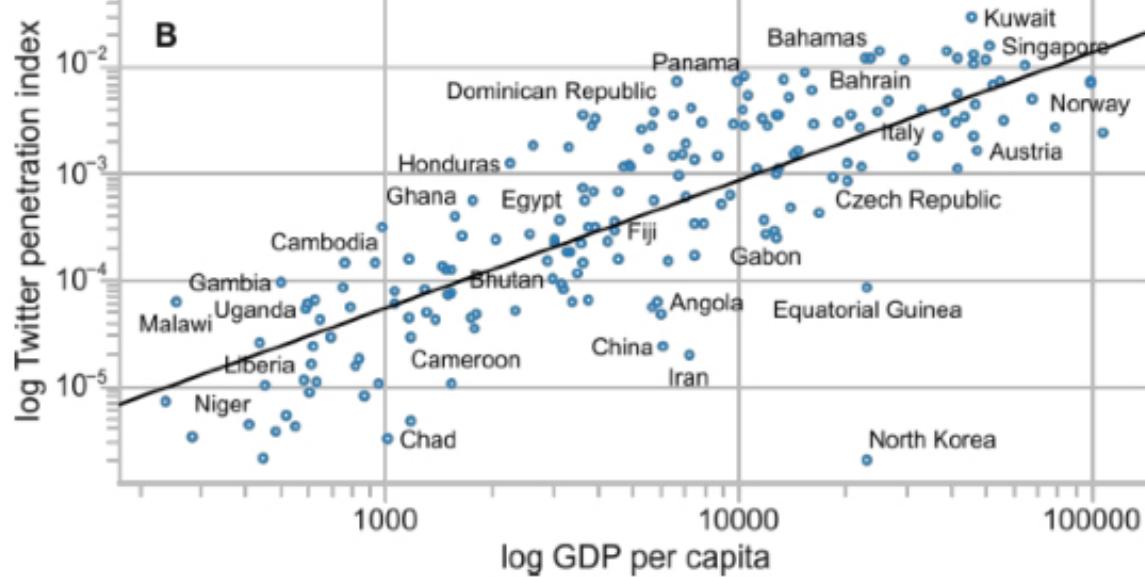
# Twitter penetration

- Is Twitter penetration related to economical development of areas?

- At country scale twitter penetration ~ GDP

Hawelka, B. et al., 2013. Geo-located Twitter as the proxy for global mobility patterns.

- At small scale is the opposite! twitter penetration ~ unemployment



# Twitter social interactions

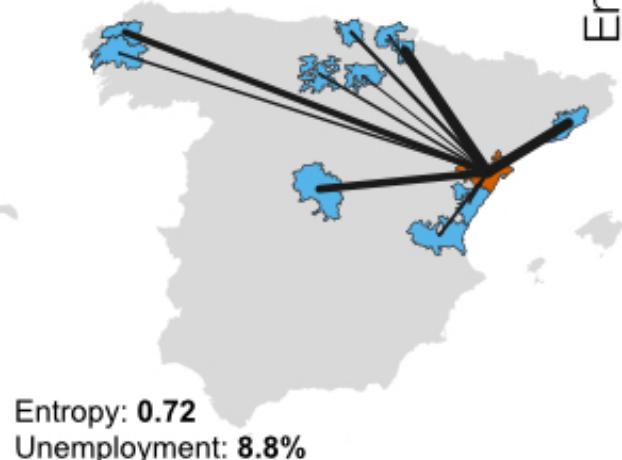
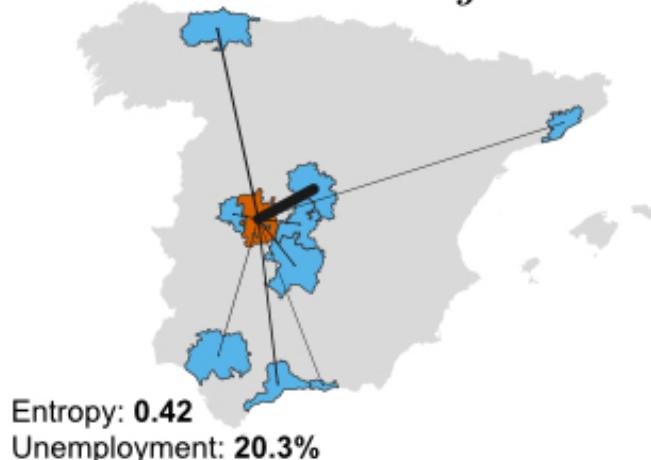
- Granovetter: diversity of interactions yields to more opportunities
- Diversity of interactions between cities is correlated with economical development  
Eagle et al, Science 2010
- We construct the graph of social interactions

$w_{ij}$  = number of @ between areas  $i$  and  $j$

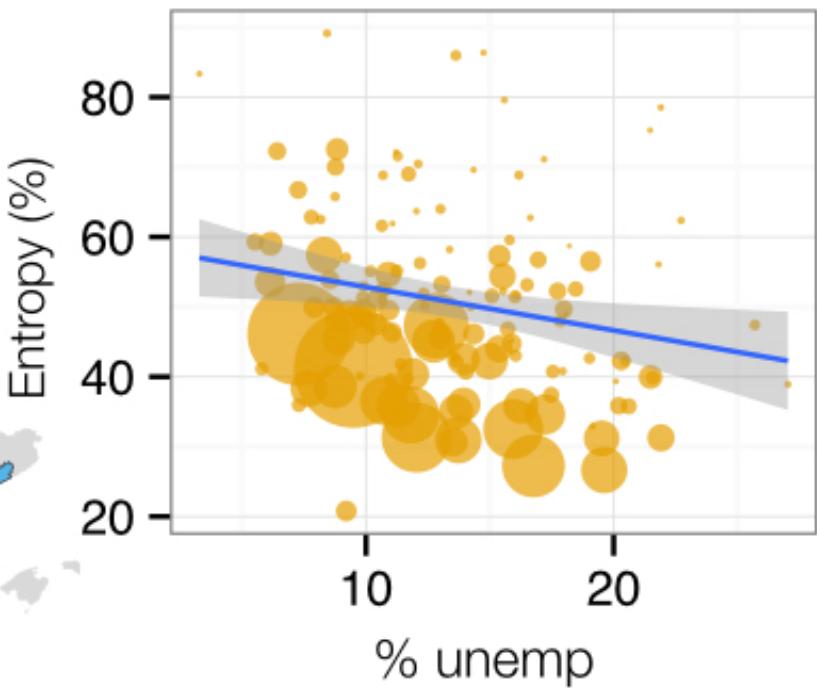
$$p_{ij} = w_{ij} / \sum_{j=1}^{k_i} w_{ij}$$

- Measure diversity with **entropy**

$$S_i = - \sum_{j=1}^{k_i} p_{ij} \log p_{ij}$$



$$\rho = -0.21[-0.37, -0.04]$$



# Twitter geographical interactions

- Diversity of geographical mobility is correlated with development

Smith, C., Mashhadi, A. & Capra, L., 2013. Ubiquitous sensing for mapping poverty in developing countries.

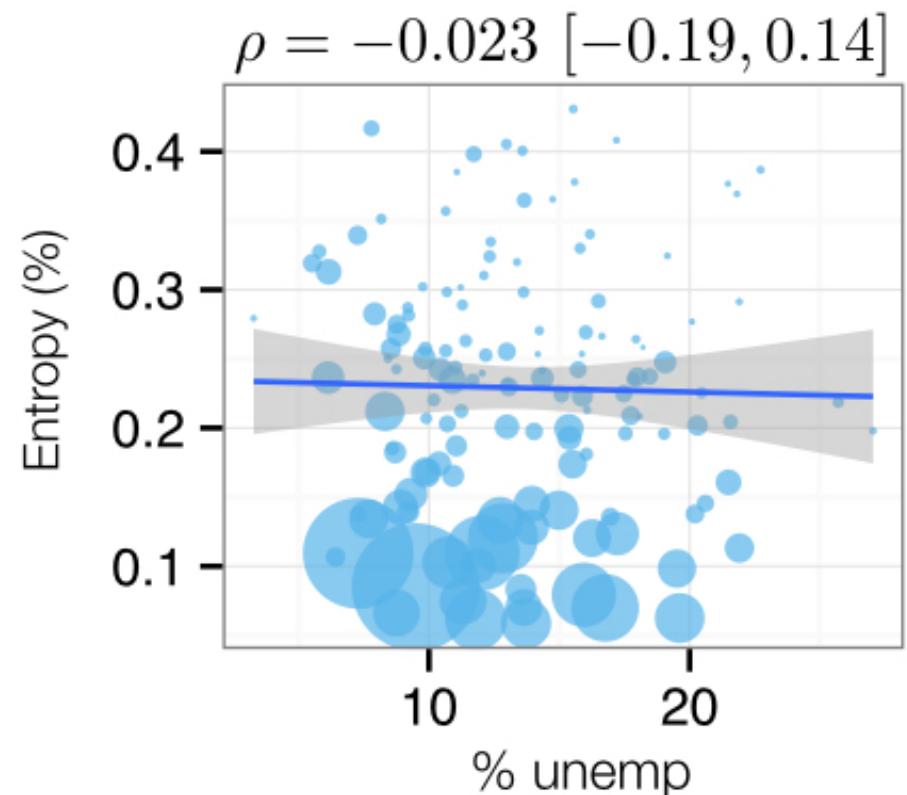
Smith, C., Quercia, D. & Capra, L., 2013. Finger on the pulse: identifying deprivation using transit flow analysis.

- We use the graph of flows

$$\tilde{p}_{ij} = T_{ij} / \sum_{j=1}^{\tilde{k}_i} T_{ij}$$

- Measure diversity with **entropy**

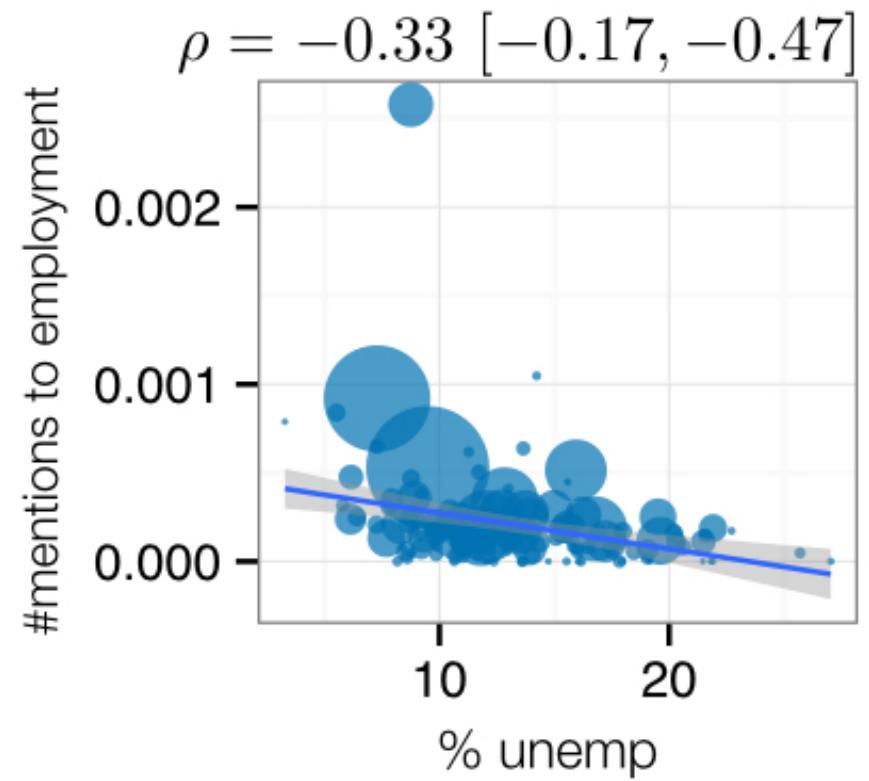
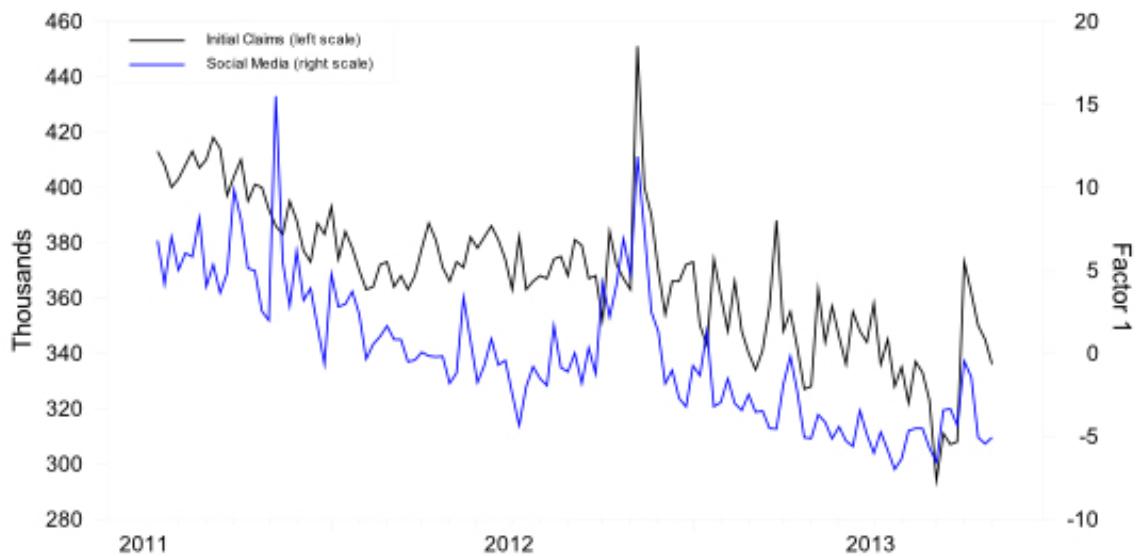
$$\tilde{S}_i = - \sum_{j=1}^{\tilde{k}_i} \tilde{p}_{ij} \log \tilde{p}_{ij}$$



# Twitter **content**

- Two different approaches
    - Classical approach: NLP applied to detect mentions to “unemployment”, “job”, “economy”, ...

Antenucci, D. et al., 2014. Using Social Media to Measure Labor Market Flows.



# Twitter content

- Our approach: NLP applied to detect **lexical complexity** (as a proxy for educational level)

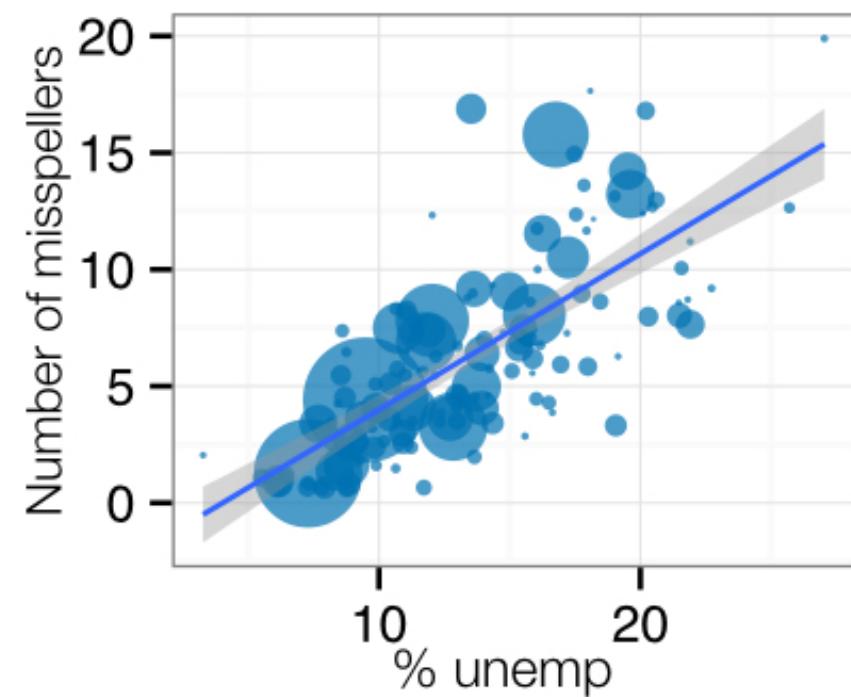
- Readability (Gunning index)

*J. Davenport et al, The Readability of Tweets and their Geographic Correlation with Education, arXiv:1401.6058, 2014*

- **Serious misspellings**

Tweet	Correct spelling
Alguien se viene <b>con migo aver</b> la vida de PI??	Alguien se viene <b>conmigo a ver</b> la vida de PI??
La quiero mucho y la <b>hecho de menos</b>	La quiero mucho y la <b>echo de menos</b>
Yo <b>llendo</b> a trabajar con este tiempo	Yo <b>yendo</b> a trabajar con este tiempo

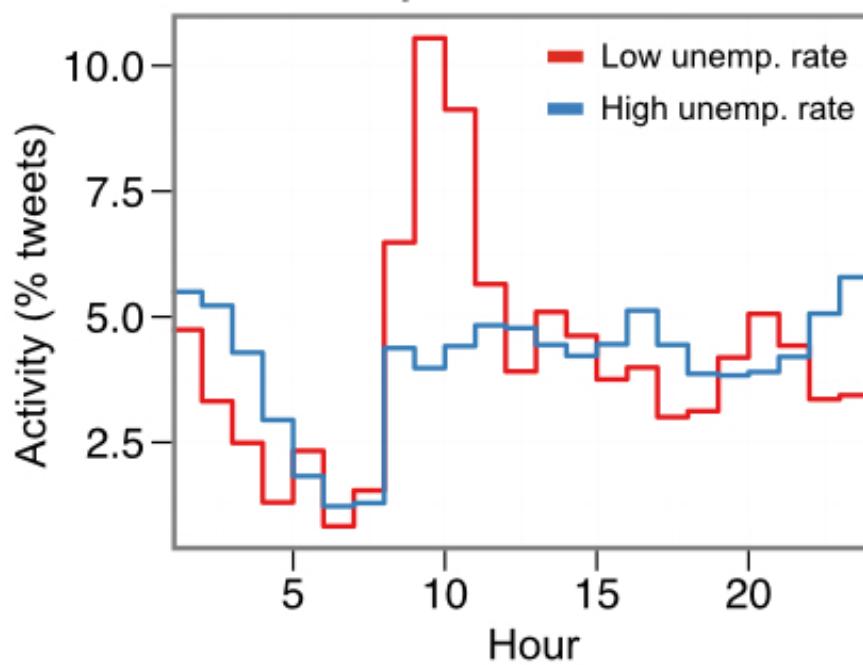
- We construct a list of more than 600 incorrect expressions of this type validated by spanish language linguistic experts.
- We do not take into account misspellings due to different Spanish accents and IM abbreviations
- We compute for each area the fraction of users that make a number of serious misspellings



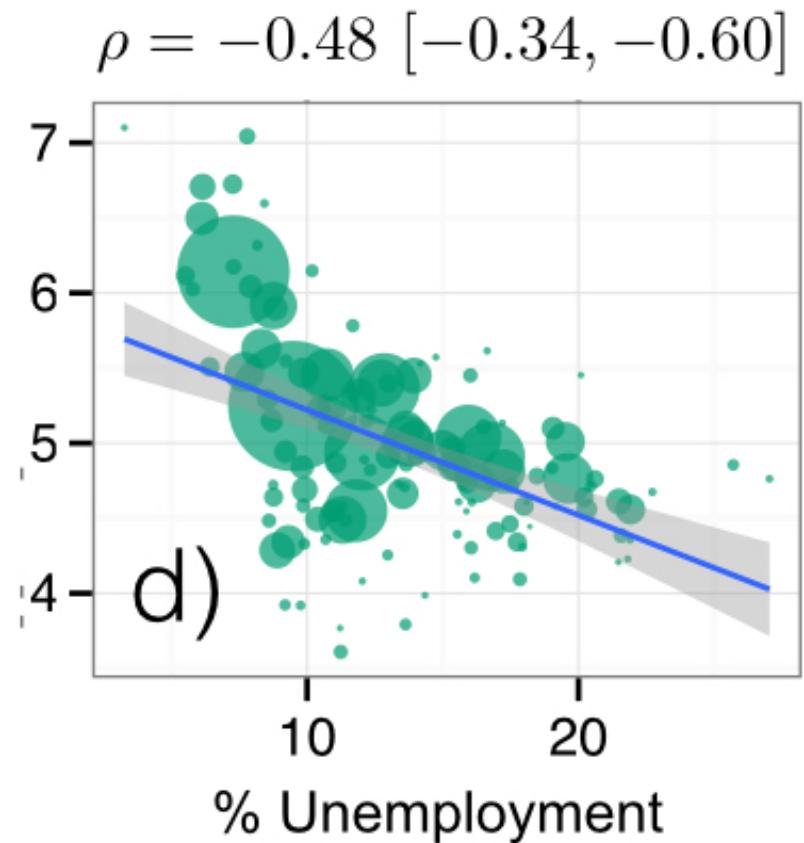
# Twitter **activity**

- Is unemployment reflected in twitter daily patterns?

Just arrived to work, mondays are too hard...  
Tweet 10:43 - 2 de jun. de 2014



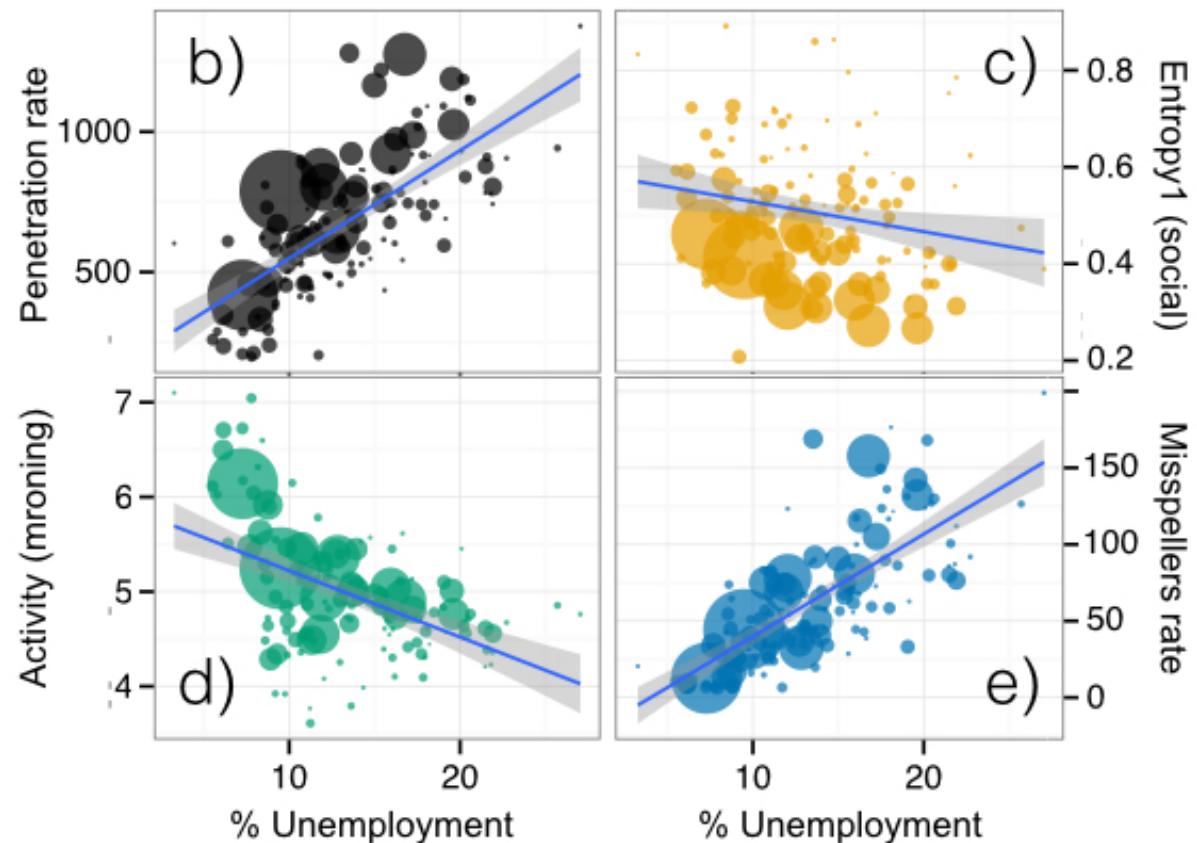
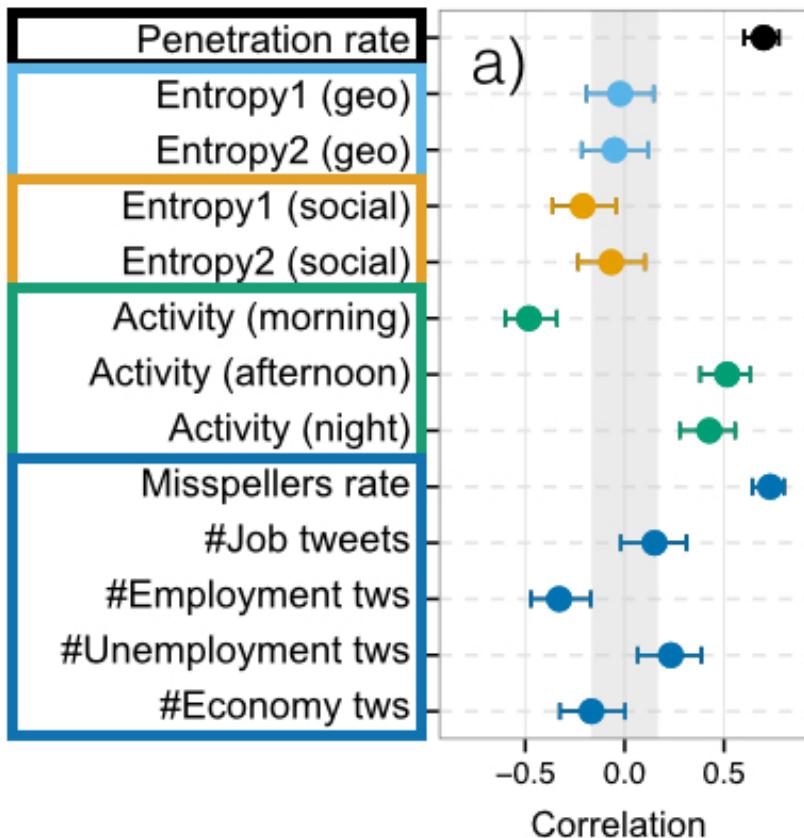
Activity (morning)



# Summary of the variables

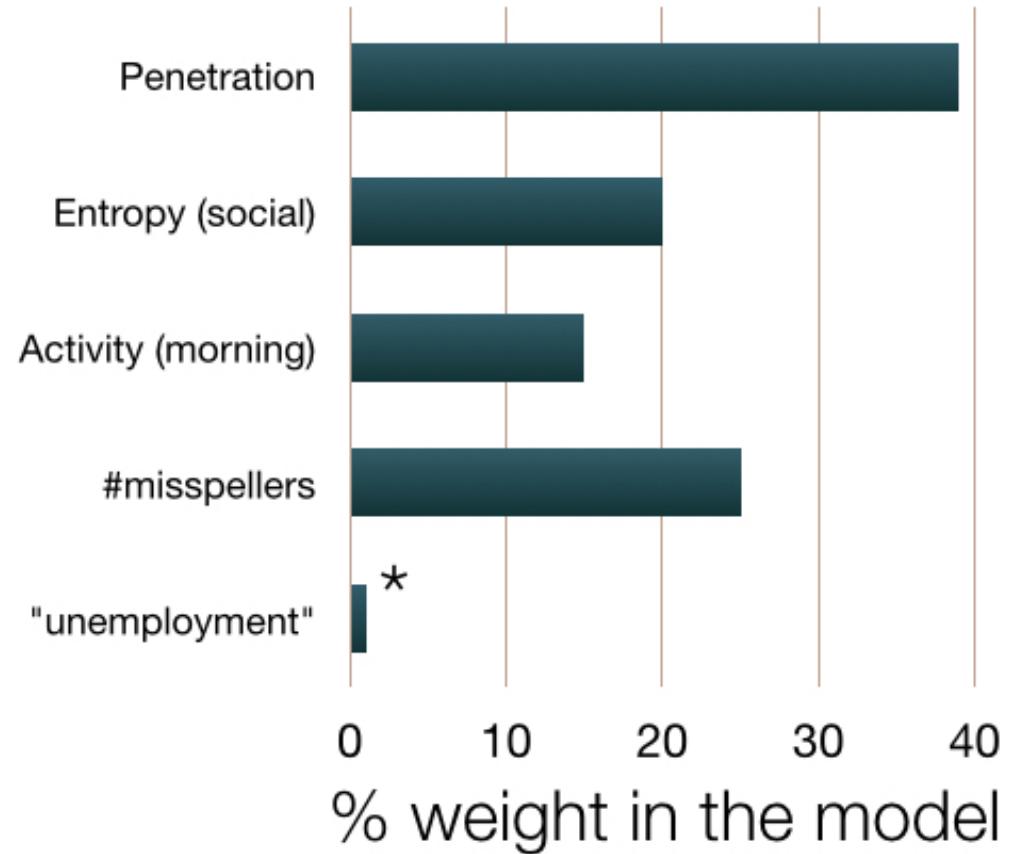
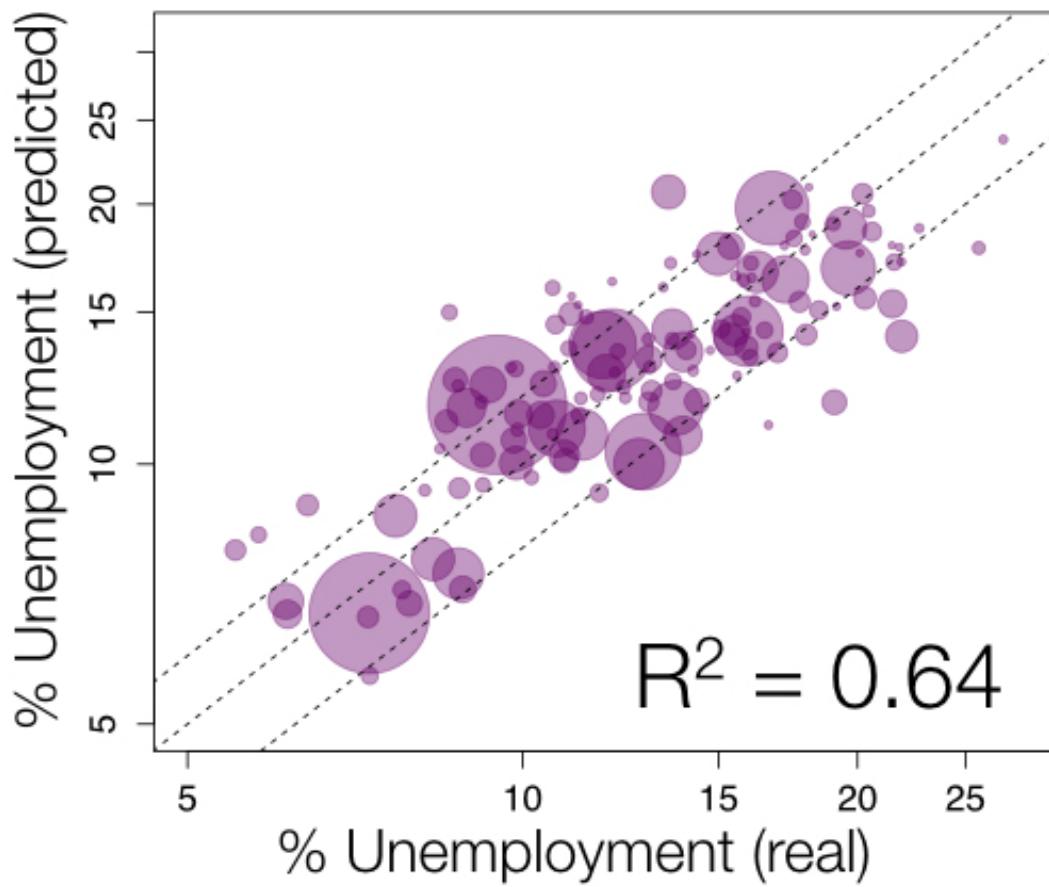
Social/geo variables have low correlation

Penetration rate/activity and content are highly correlated with unemployment



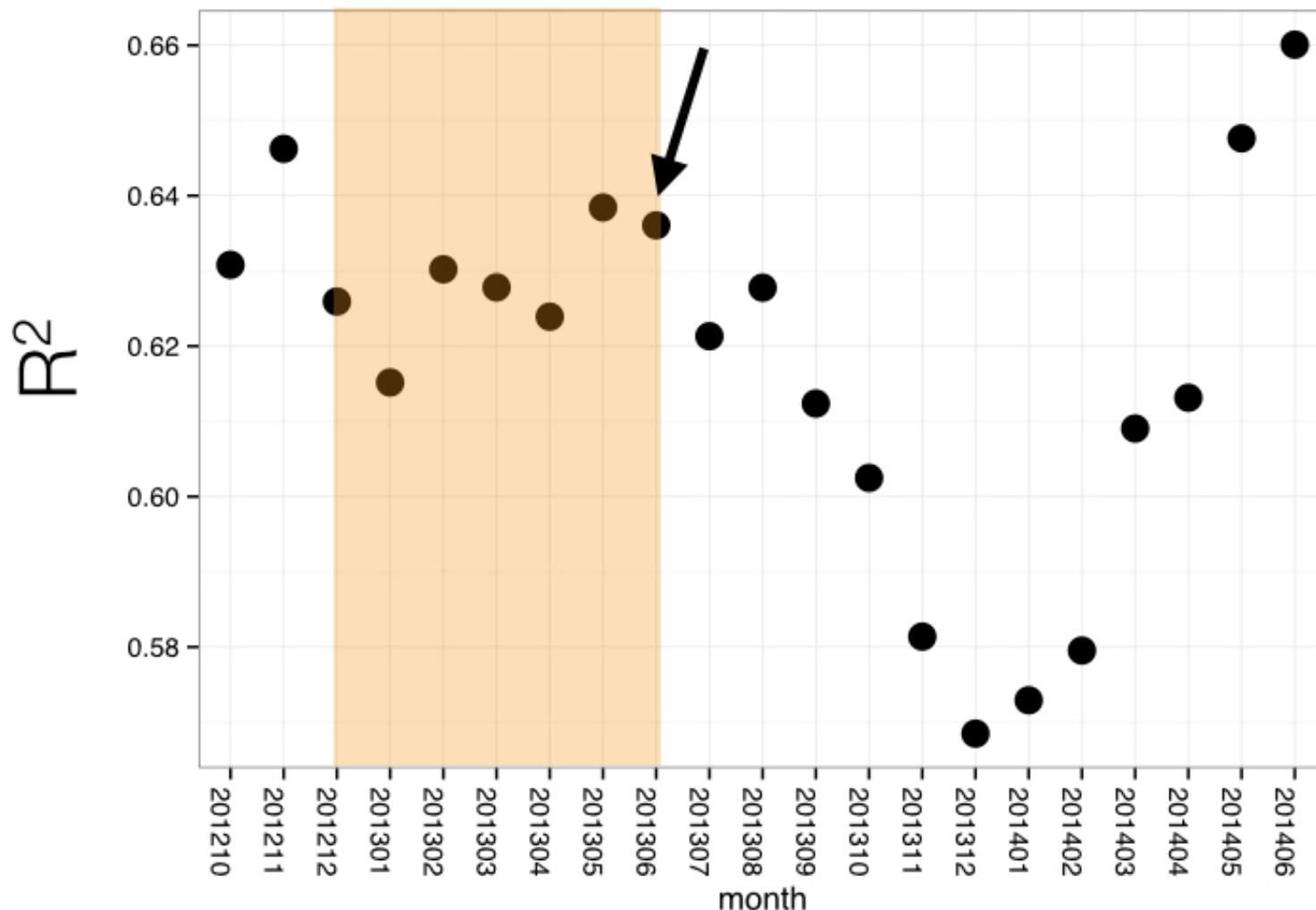
# Explanatory power of Twitter variables

- Simple linear regression



# Explanatory power as a function of time

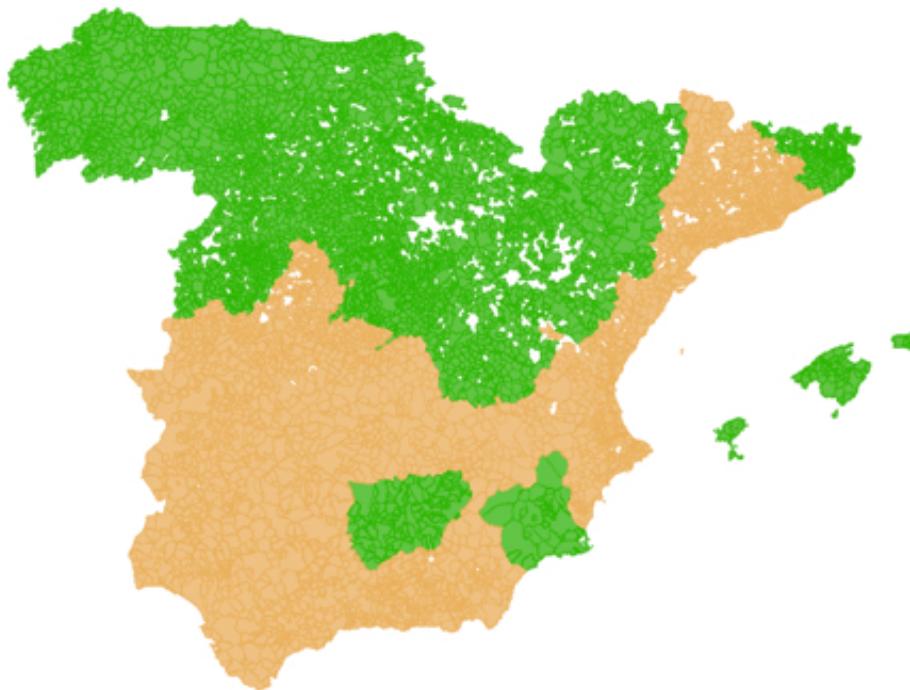
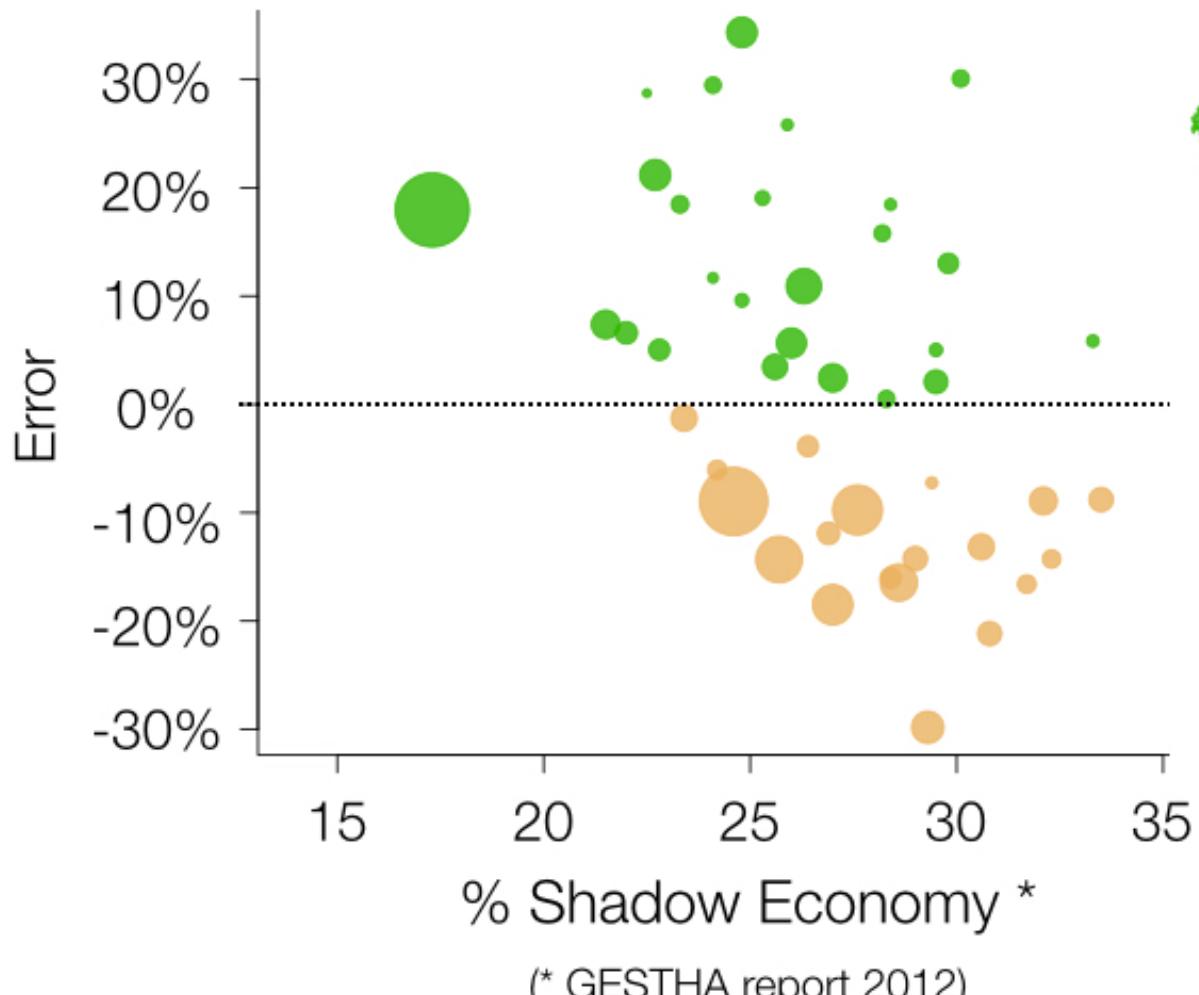
- At what time of the year the model has more explanatory power?



# Are we really wrong?

---

Model Error = Model[variables] - Official unemployment



# Summary

---

# Summary

---

- **Economical development -> Behavior -> Social Media**

# Summary

---

- **Economical development -> Behavior -> Social Media**
- Can Twitter be used to infer economical development? YES
  - Areas with different behavior in Twitter show different levels of unemployment
  - Applications to marketing, media, planning
  - Twitter is everywhere! and the API is free! Prove us wrong, please!

# Summary

---

- **Economical development -> Behavior -> Social Media**
- Can Twitter be used to infer economical development? YES
  - Areas with different behavior in Twitter show different levels of unemployment
  - Applications to marketing, media, planning
  - Twitter is everywhere! and the API is free! Prove us wrong, please!
- Can Twitter variables explain the unemployment per area? YES with  $R^2 = 0.64$ 
  - Activity, penetration and content account for 80% of the variance explained
  - Diversity on geographical and social interaction amount only for the 20%

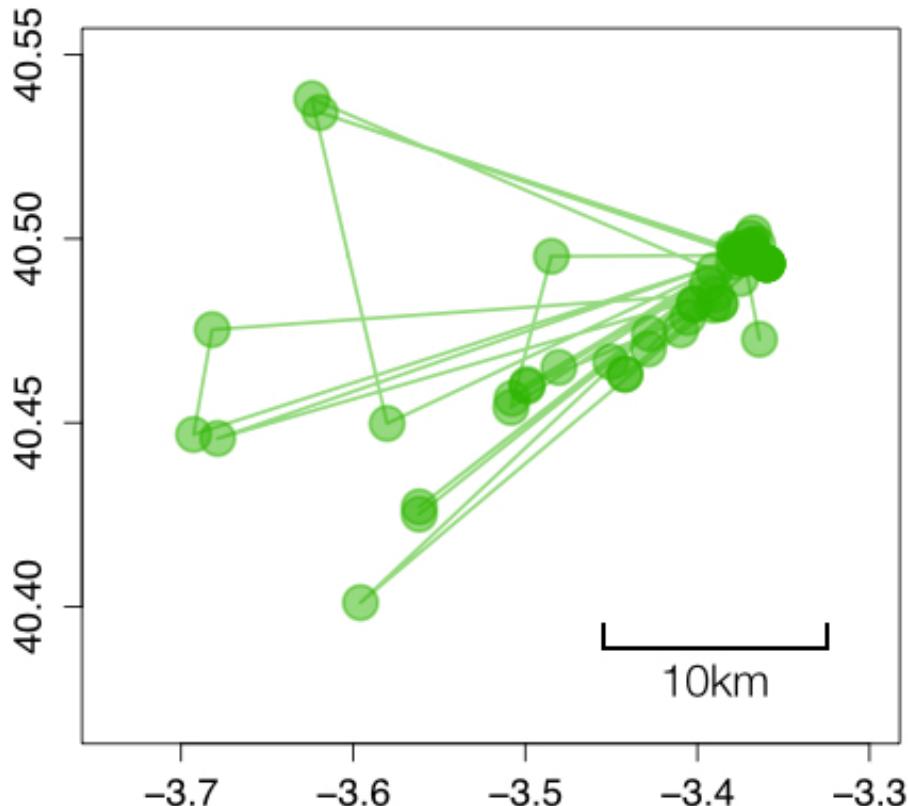
# Summary

---

- **Economical development -> Behavior -> Social Media**
- Can Twitter be used to infer economical development? YES
  - Areas with different behavior in Twitter show different levels of unemployment
  - Applications to marketing, media, planning
  - Twitter is everywhere! and the API is free! Prove us wrong, please!
- Can Twitter variables explain the unemployment per area? YES with  $R^2 = 0.64$ 
  - Activity, penetration and content account for 80% of the variance explained
  - Diversity on geographical and social interaction amount only for the 20%
- **What we are doing now:**
  - Can we use the model to forecast future unemployment? Or now-casting?
  - What are the behavioral changes behind socio-economical changes?
  - Use the model in under-developed countries.

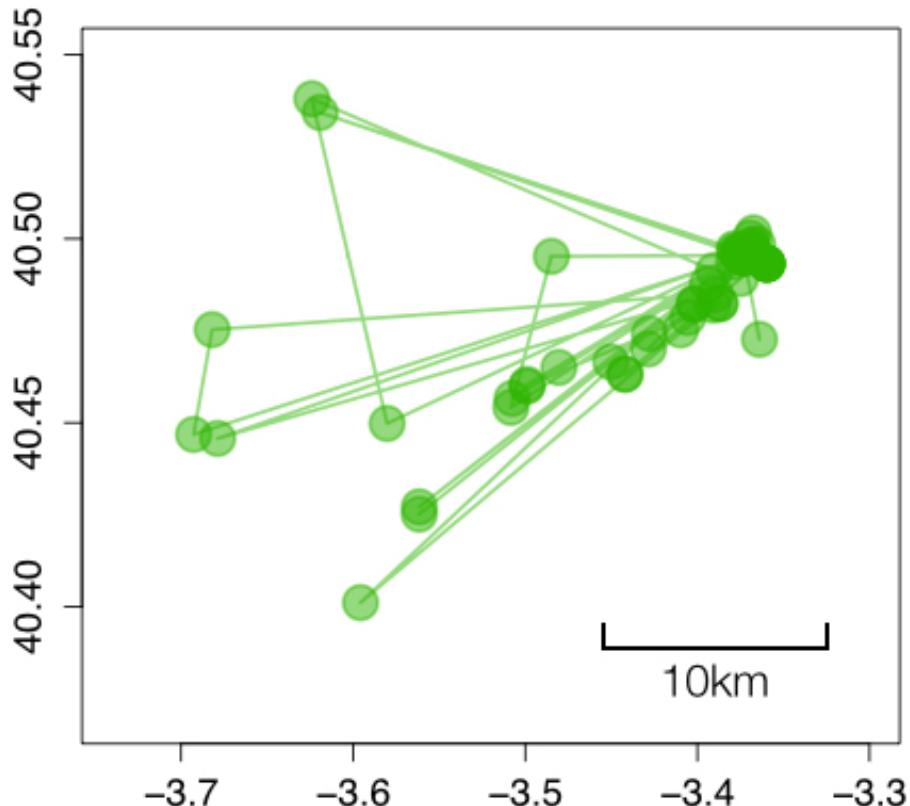
# Behavioral changes behind socio-economical changes

(Geolocalized tweets)

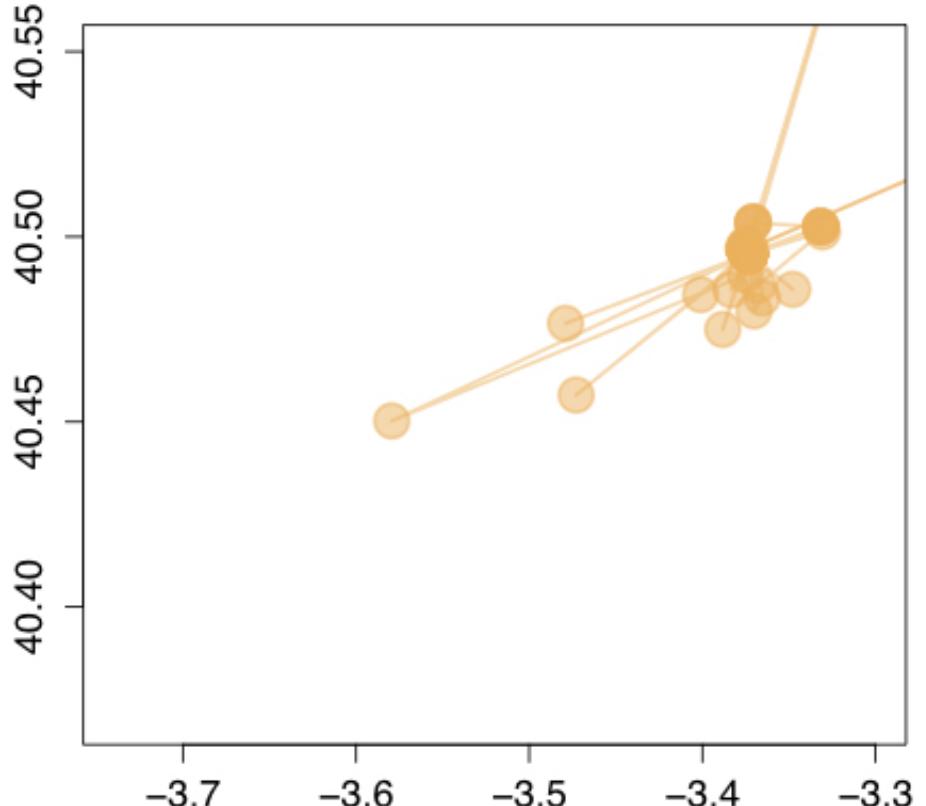


# Behavioral changes behind socio-economical changes

(Geolocalized tweets)

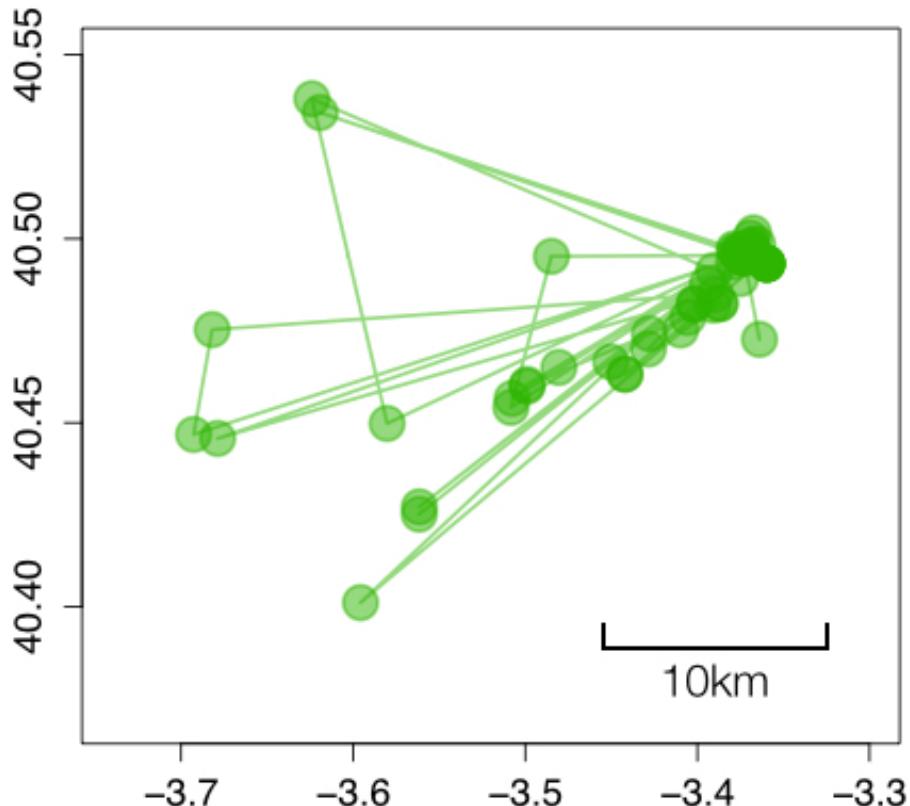


**Working**

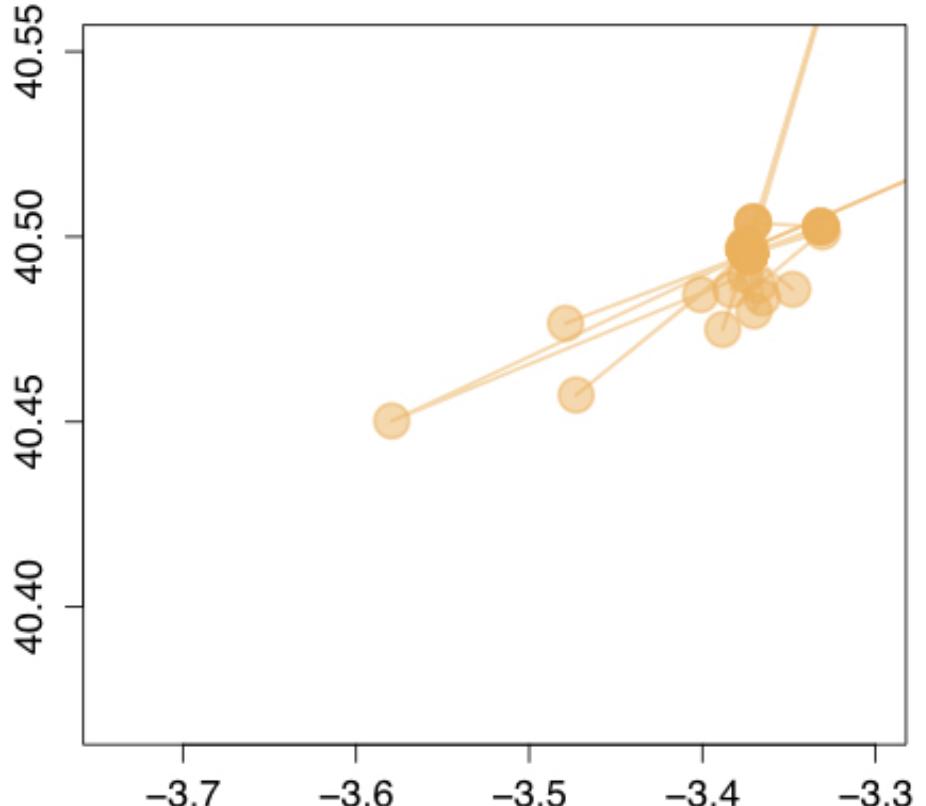


# Behavioral changes behind socio-economical changes

(Geolocalized tweets)



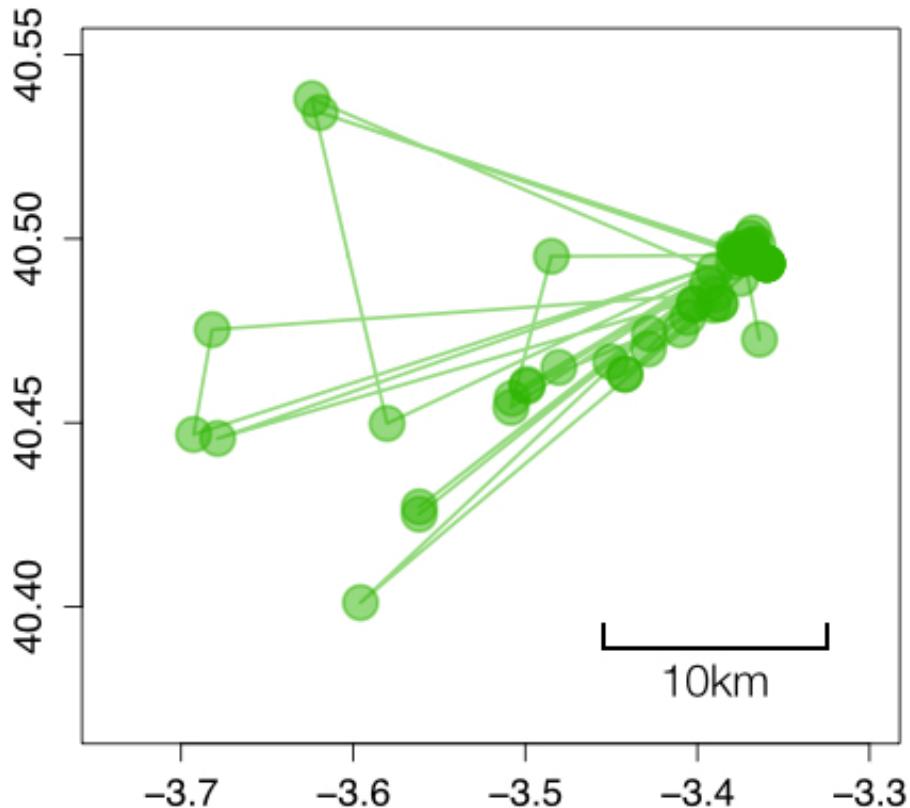
**Working**



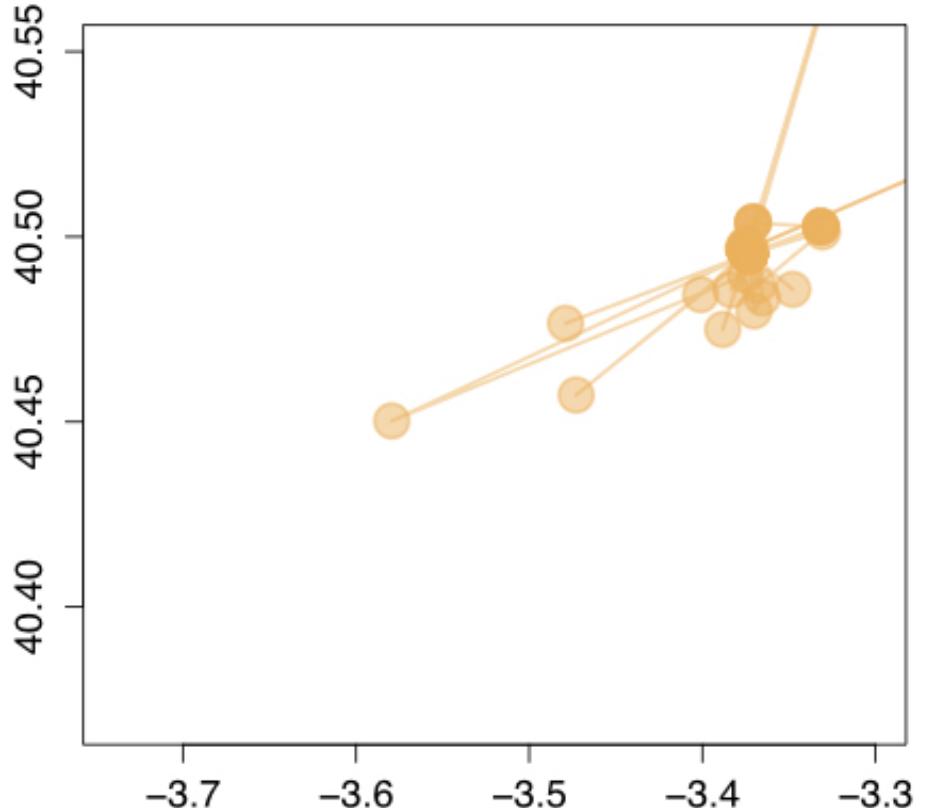
**Unemployed**

# Behavioral changes behind socio-economical changes

(Geolocalized tweets)



**Working**



**Unemployed**

Less geographical mobility, more probability to be unemployed