

# BIGDATA INEGI



## Use of data from social networks to obtain statistical and geographical information

Global Conference on Big Data for Official Statistics  
October 2015, Abu Dhabi, UAE



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# Purpose of the presentation

Sharing Mexican experiences about the use of tweets to obtain statistical information by the National Institute of Statistics and Geography of Mexico, INEGI.

Pilot Test



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# CONTEXT



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# National Institute of Statistics and Geography of Mexico, INEGI

INEGI is a Constitutionally **Autonomous Entity** with:

- **Statistical and geographical** responsibilities
- Around 20,000 employees
- **10 Regional Offices** and **34 State Offices** all around the country
- **5 General Directorates** which generate and integrate statistical data: economic, social, demographic, government, public safety and justice

INEGI is responsible for Mexico's **National System of Statistical and Geographical Information**



# Emergence and Evolution of new information sources

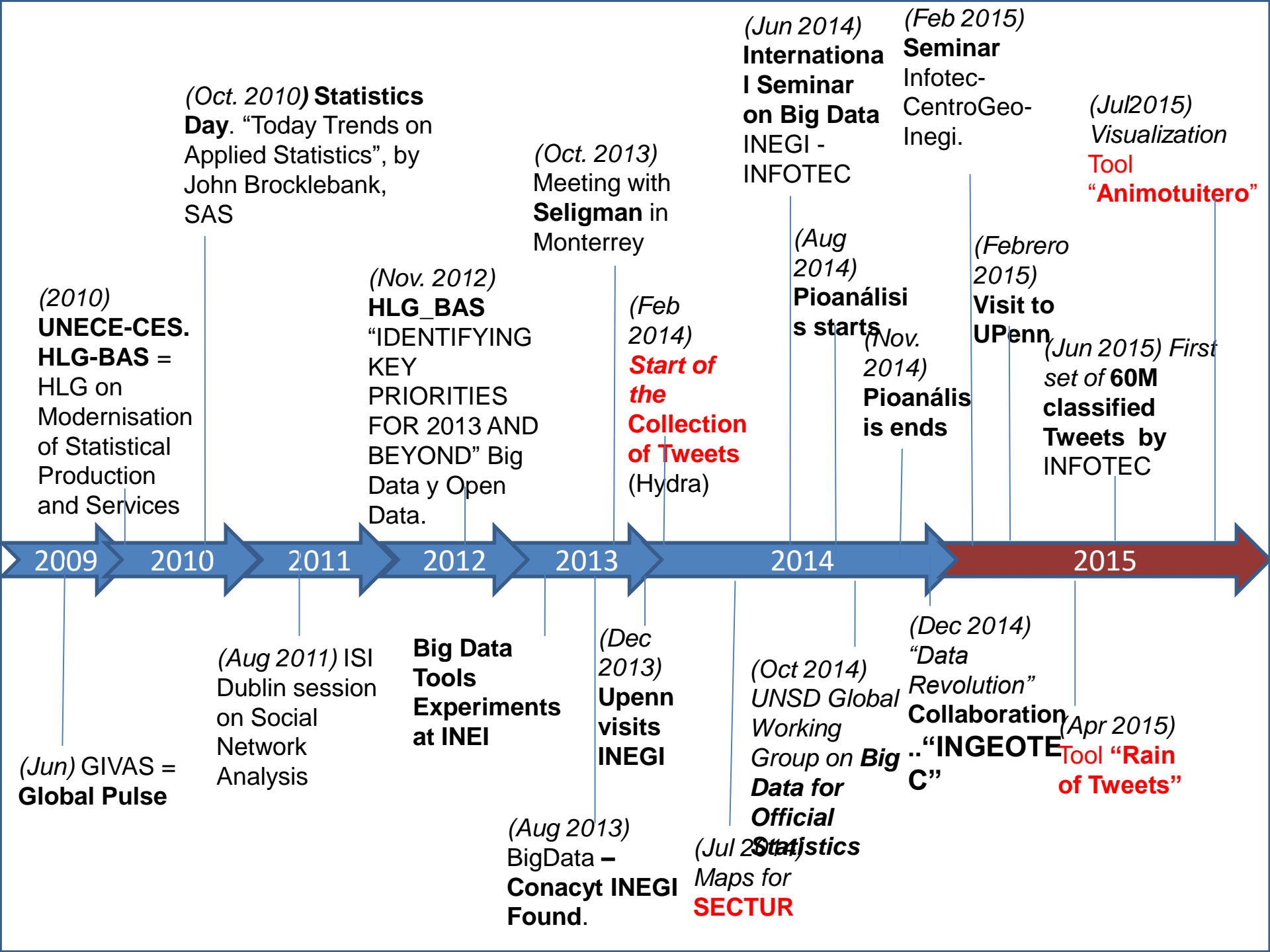


INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# Big Data Sources

- Smart meters and sensors
  - Traffic cameras, GPS devices, price scanners, power monitors, smart watches, smart phones, etc.
- Social Interactions
  - Talks and publications on social networks like Twitter, Facebook, FourSquare, etc.
- Business Transactions
  - Movements of credit cards, electronic cash registers, cell phone records, etc.
- Electronic files
  - Documents which are available in electronic formats such as PDF files, websites, videos, audio, digital media broadcasting
- Broadcast media
  - Digital video and audio produced on real time







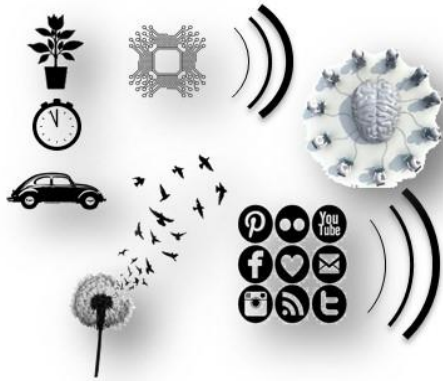
# Governing Board Commitment





# Technological Landscape

## Internet of things, people and ideas

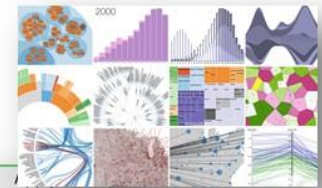
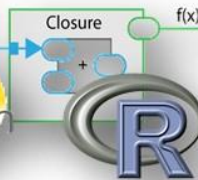
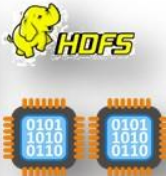
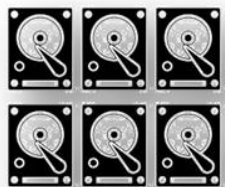


## Advanced statistics, mathematics and data Science

Signal processing  
Probability models  
Machine learning  
Statistical learning  
Data mining  
Database  
Data engineering  
Pattern recognition  
Learning patterns  
Predictive analytics

Uncertainty modeling  
Data warehousing  
Data compression  
computer  
Programming  
High-performance  
computing  
Geolocation  
Geo-referencing  
...

## Business Knowledge (experts)



**IT Infrastructure: Robust Computing and Communications, Specialized software tools for processing, analysis, visualization, etc.**

# Institutional Cooperation

- National
- International



Innovación con propósito de vida.



UNITED NATIONS  
ECONOMIC COMMISSION FOR EUROPE



United Nations  
Statistics Division



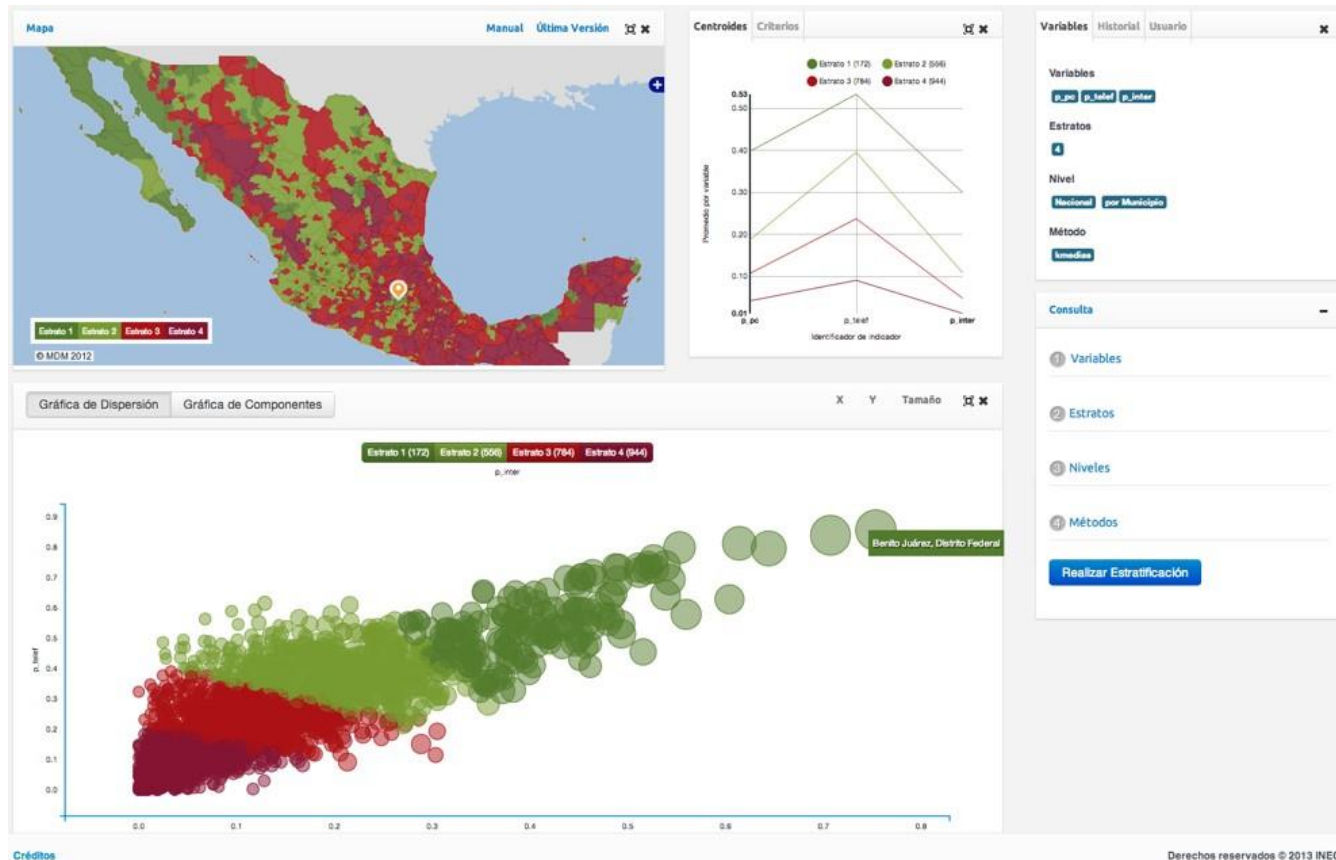
INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# TWEETS COLLECTION



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# Stratification: First Efforts

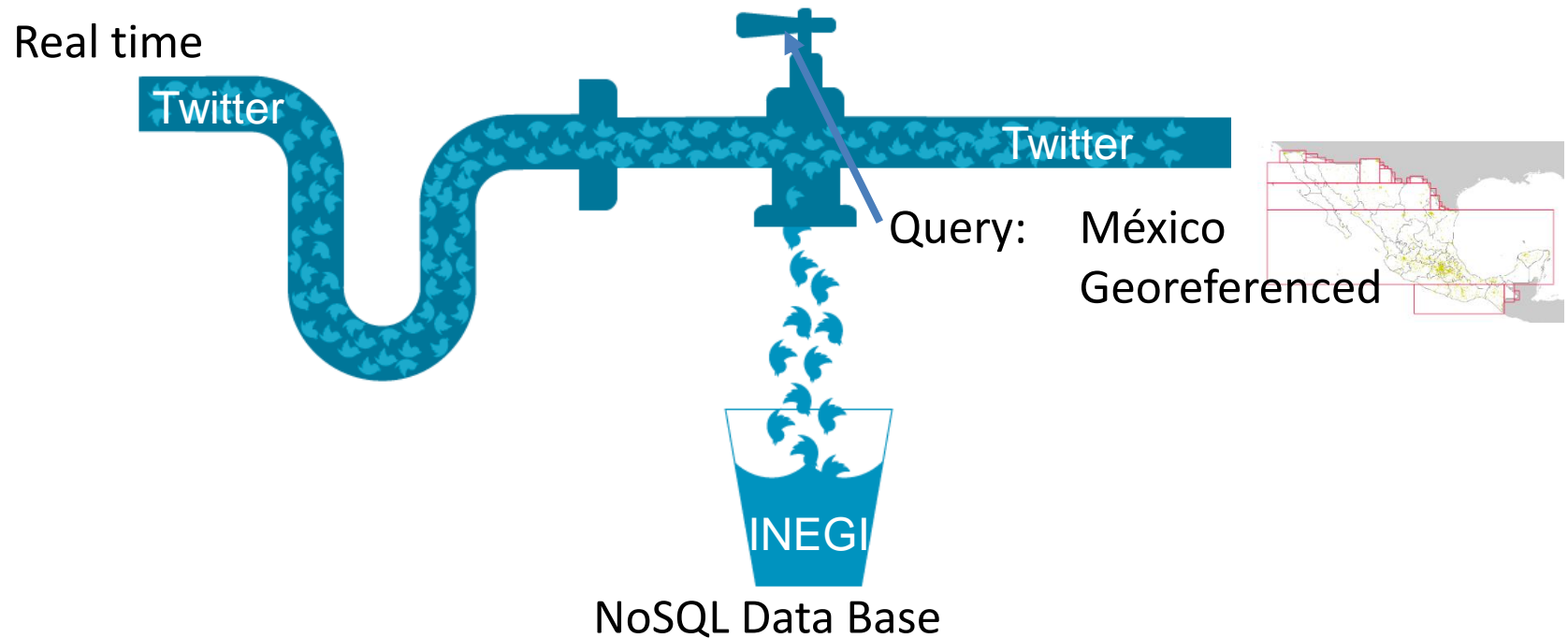


[www.inegi.org.mx/est/contenidos/Proyectos/estratificador/](http://www.inegi.org.mx/est/contenidos/Proyectos/estratificador/)



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# Twitter as a data source



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

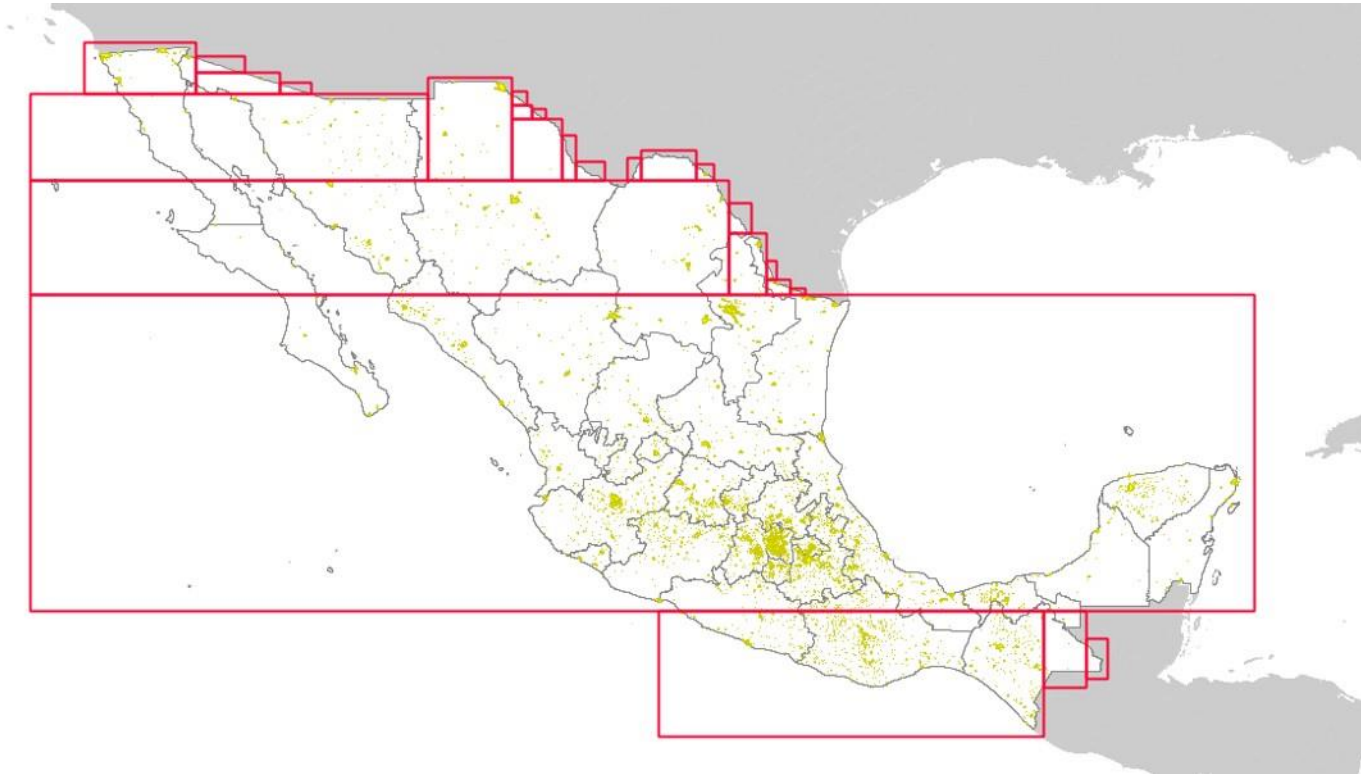
# Why Twitter?

- Readily available
- Up to 1% of global tweets at no cost
- Around 12 M accounts in Mexico
- Geolocated tweets by 700 thousand accounts
- 150 M plus tweets downloaded since January 2014
- Even though its drawbacks: Not documented, not supported by “traditional” statistics methodologies



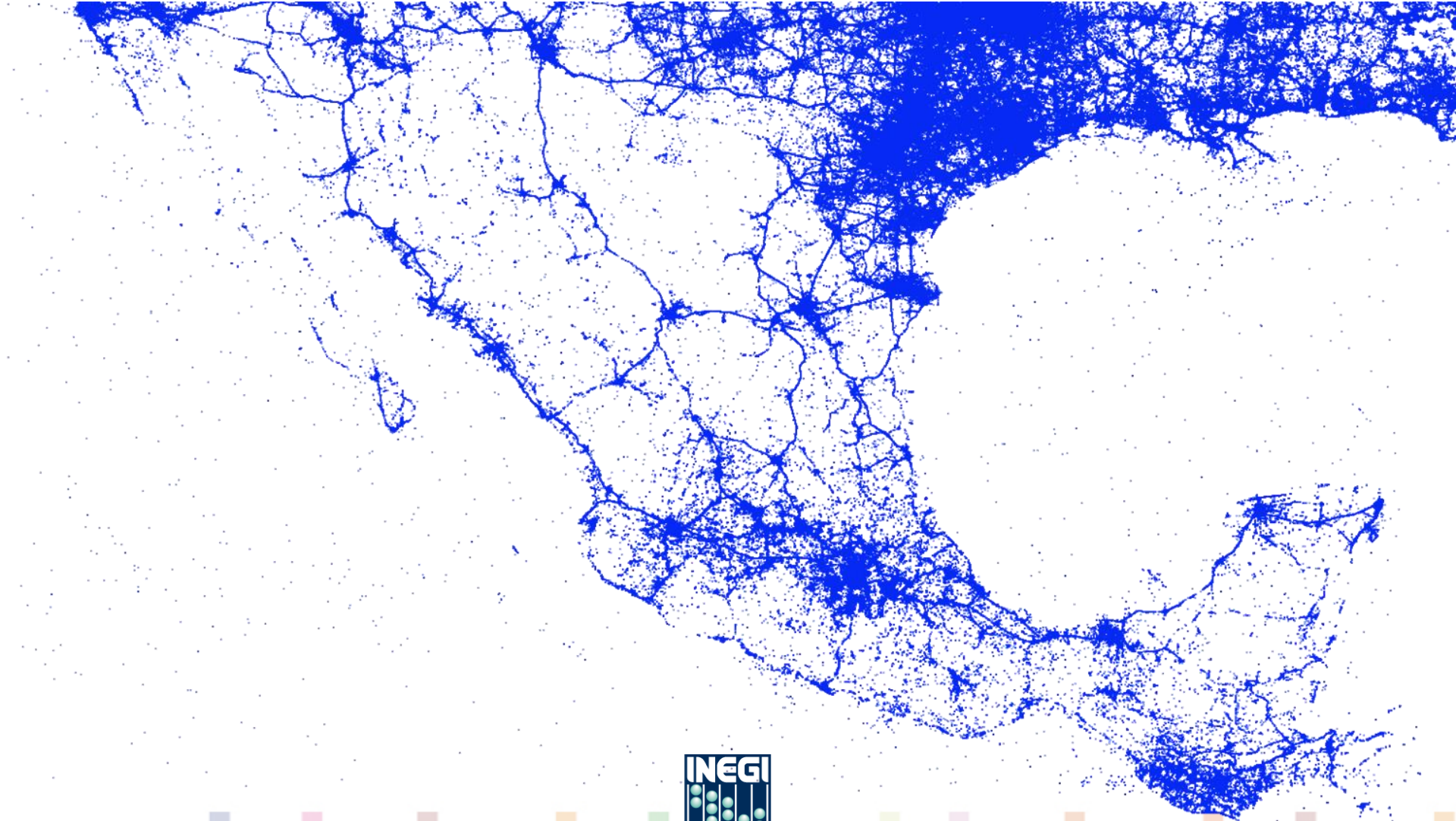


# Polygon for the collection of tweets





# 150 Millions of Tweets, August 2015

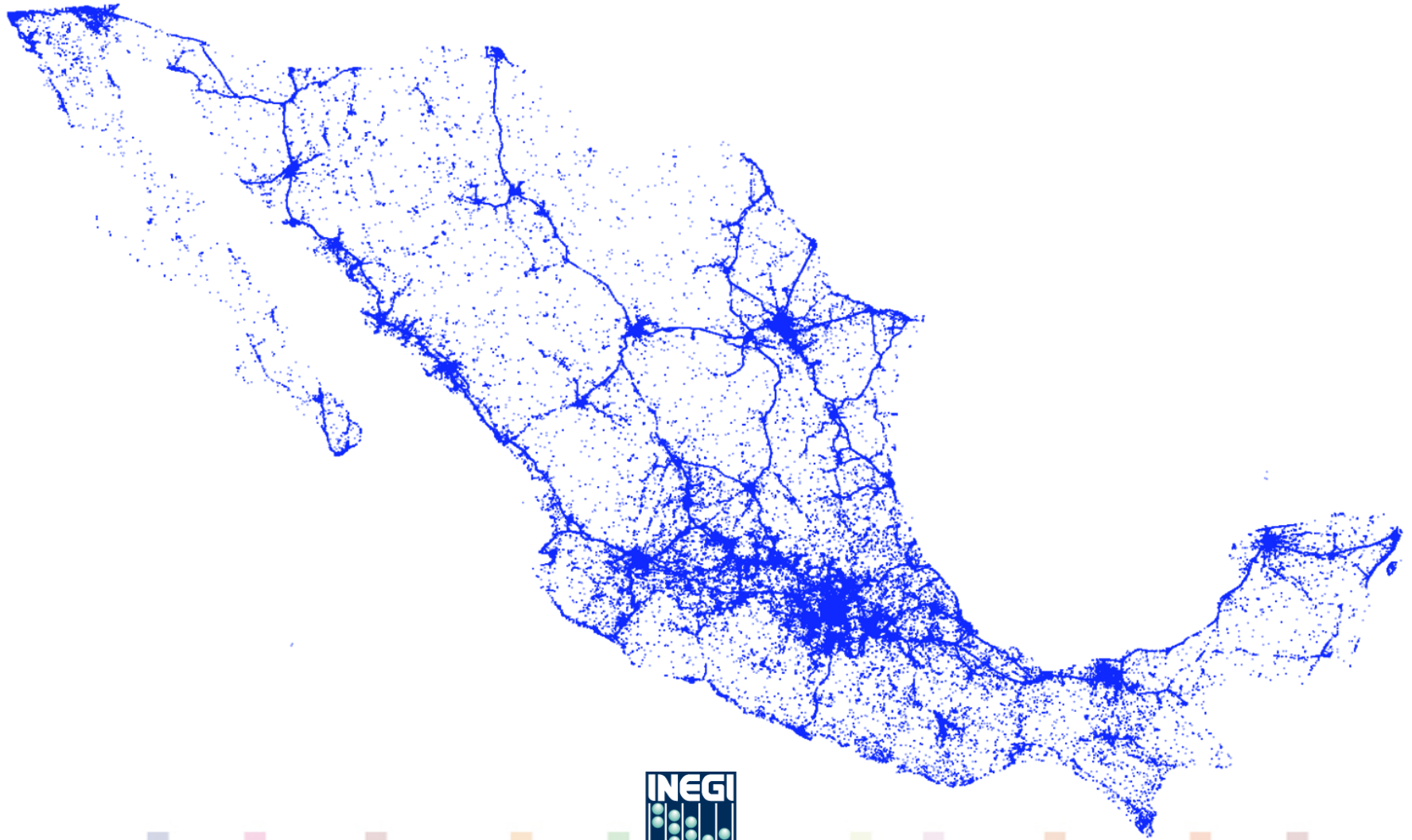


INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA



# 70+ Millions of Tweets

## August 2015



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# MOBILITY

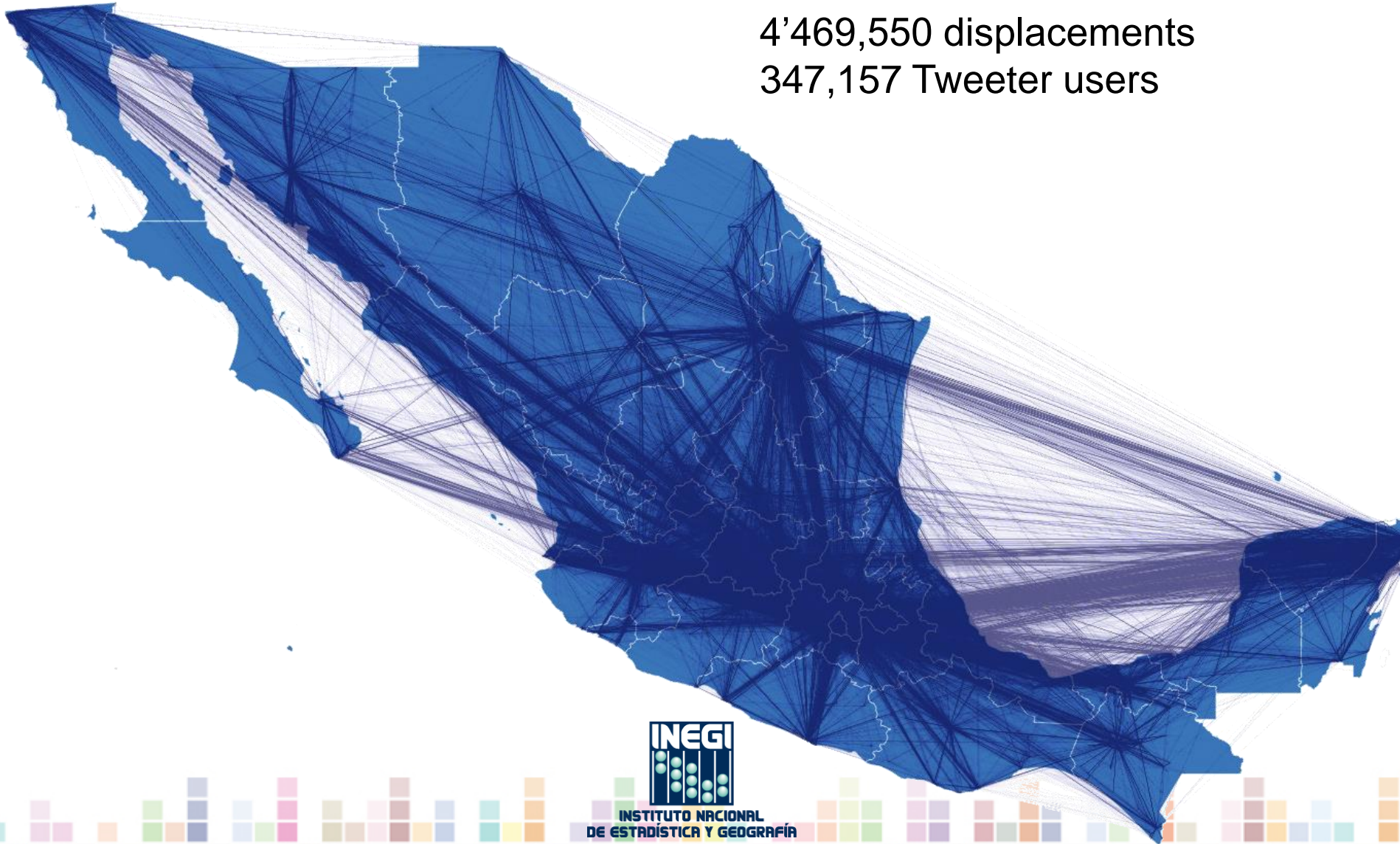


INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA



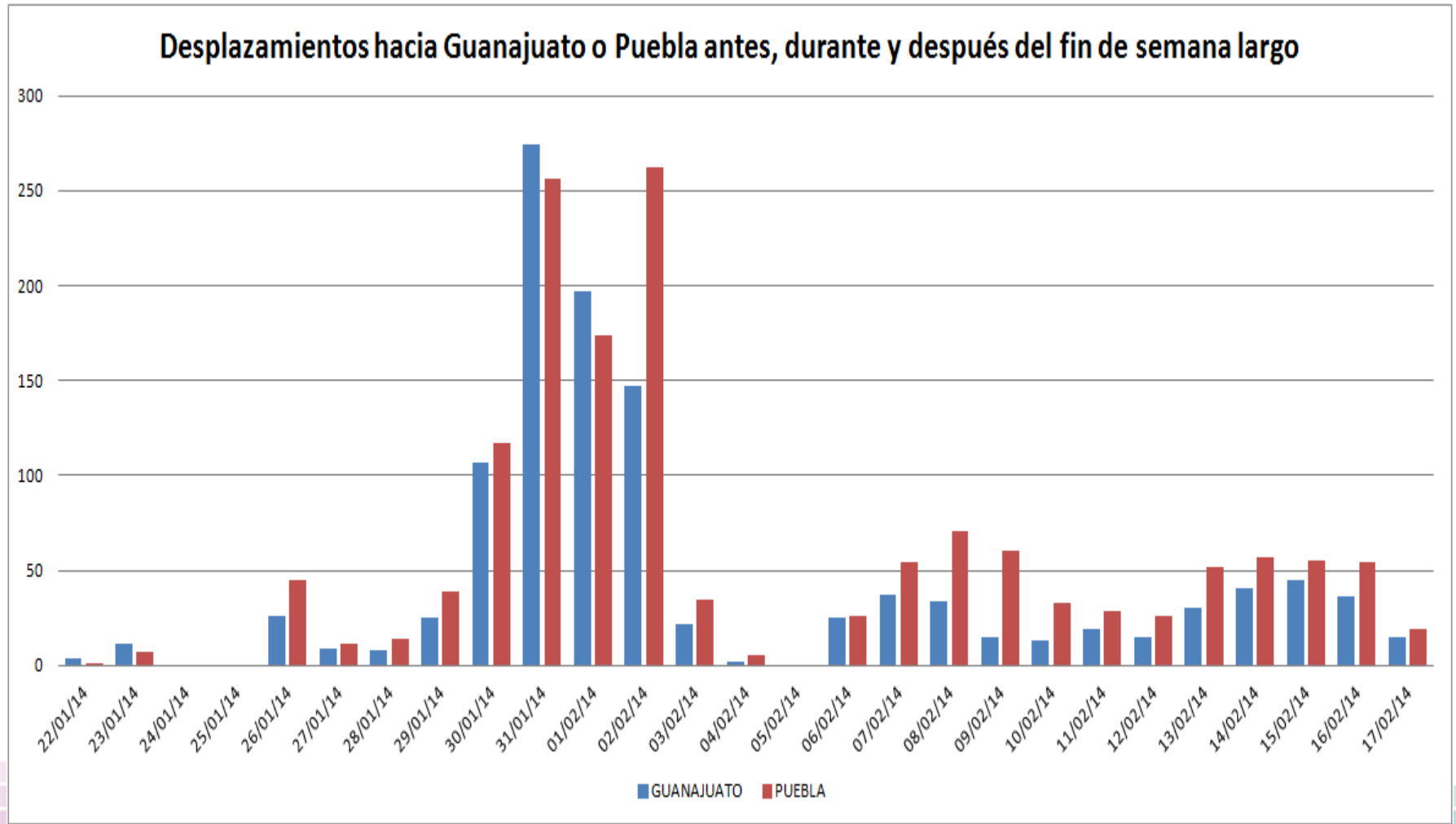
# Tweeter Users Mobility

4'469,550 displacements  
347,157 Tweeter users

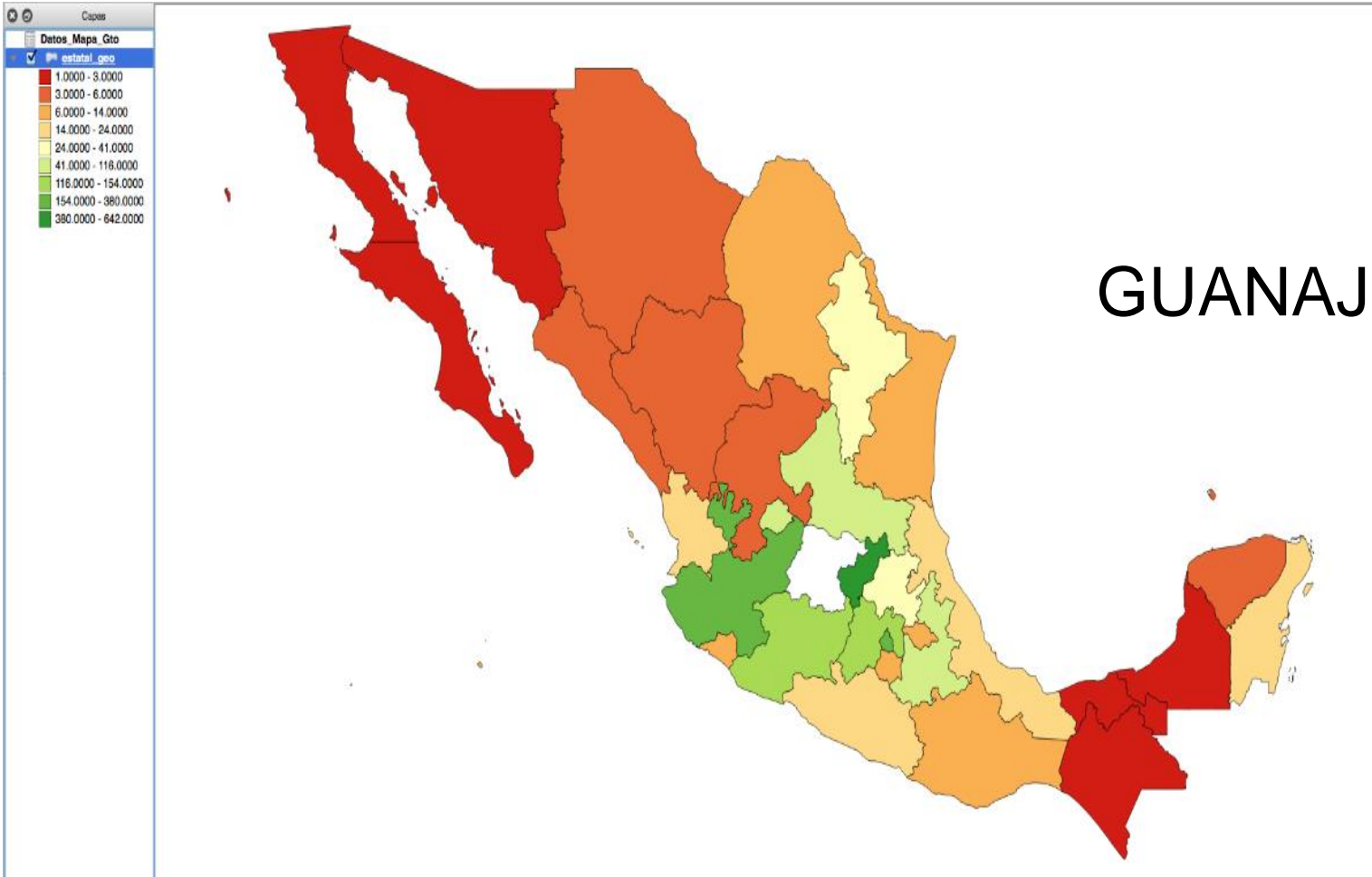


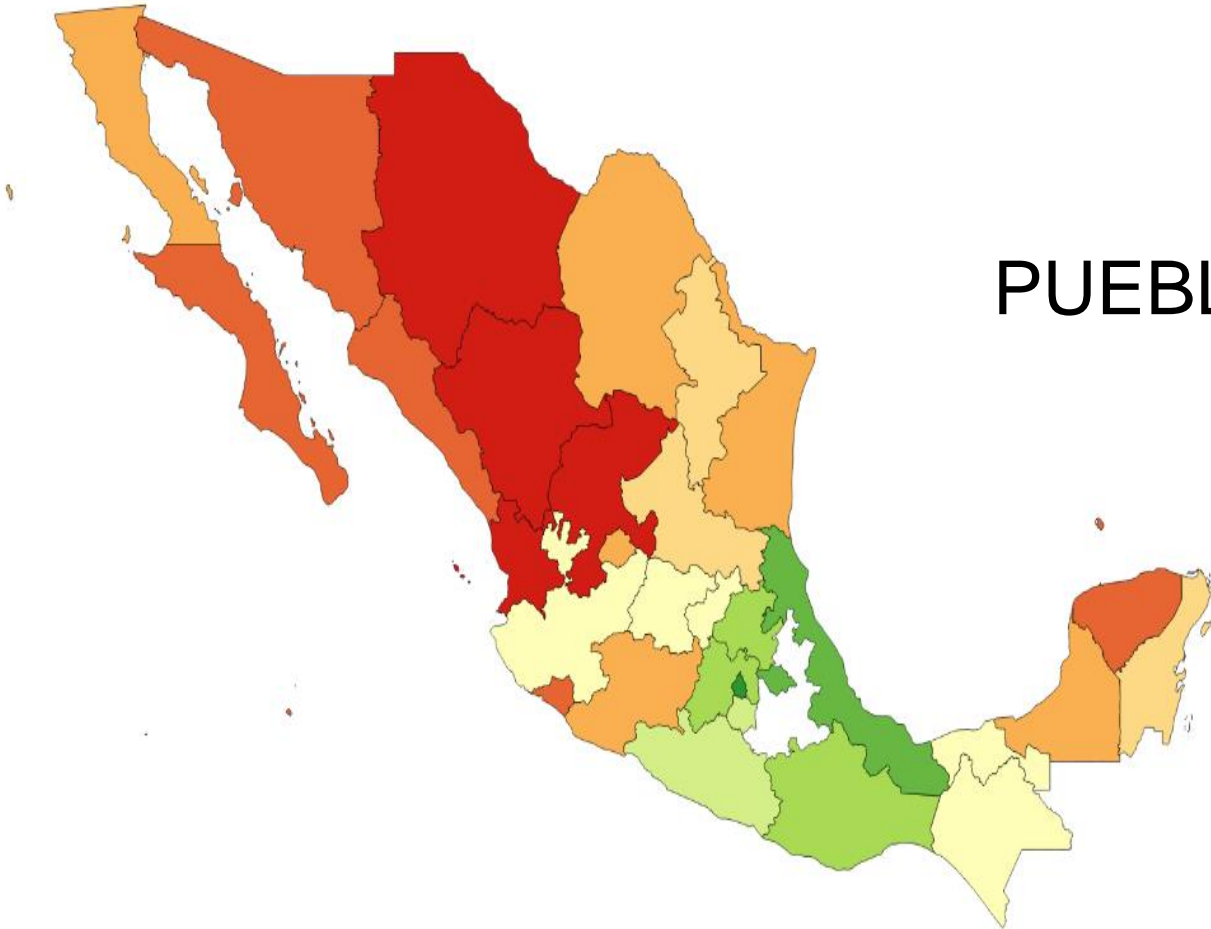
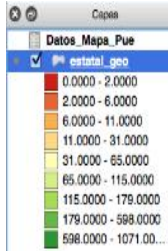
INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# Tweets behavior on a Long-Weekend



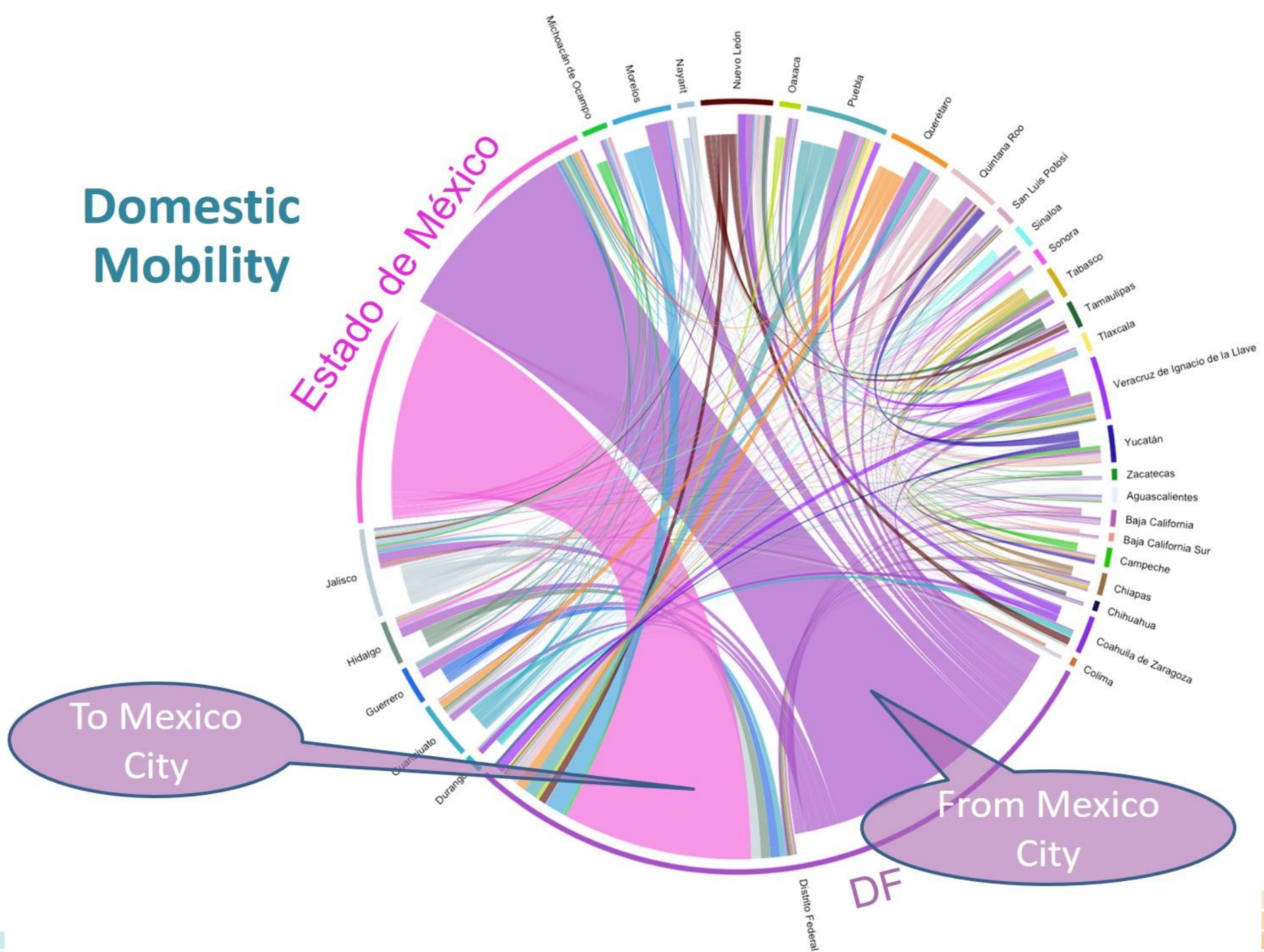
# Use of Twitter Data for Tourism Studies





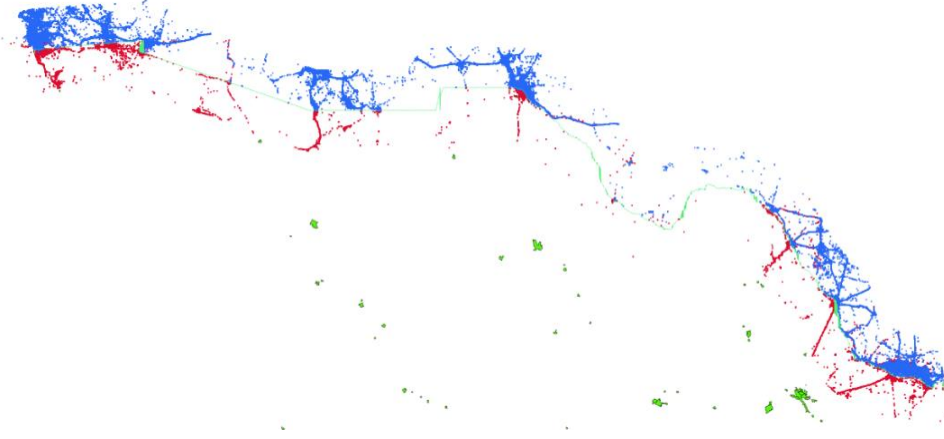


# Domestic Mobility

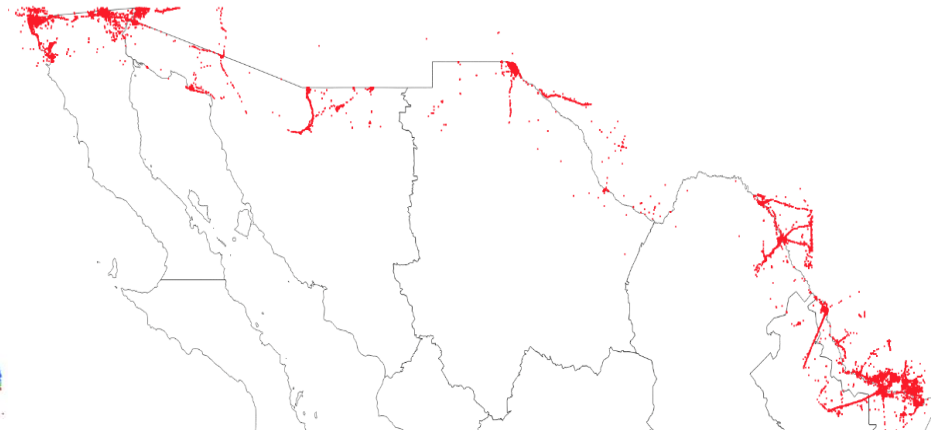


# Border Mobility

Research for development of analytical method to measure trans-border mobility through Geo-referenced tweets.



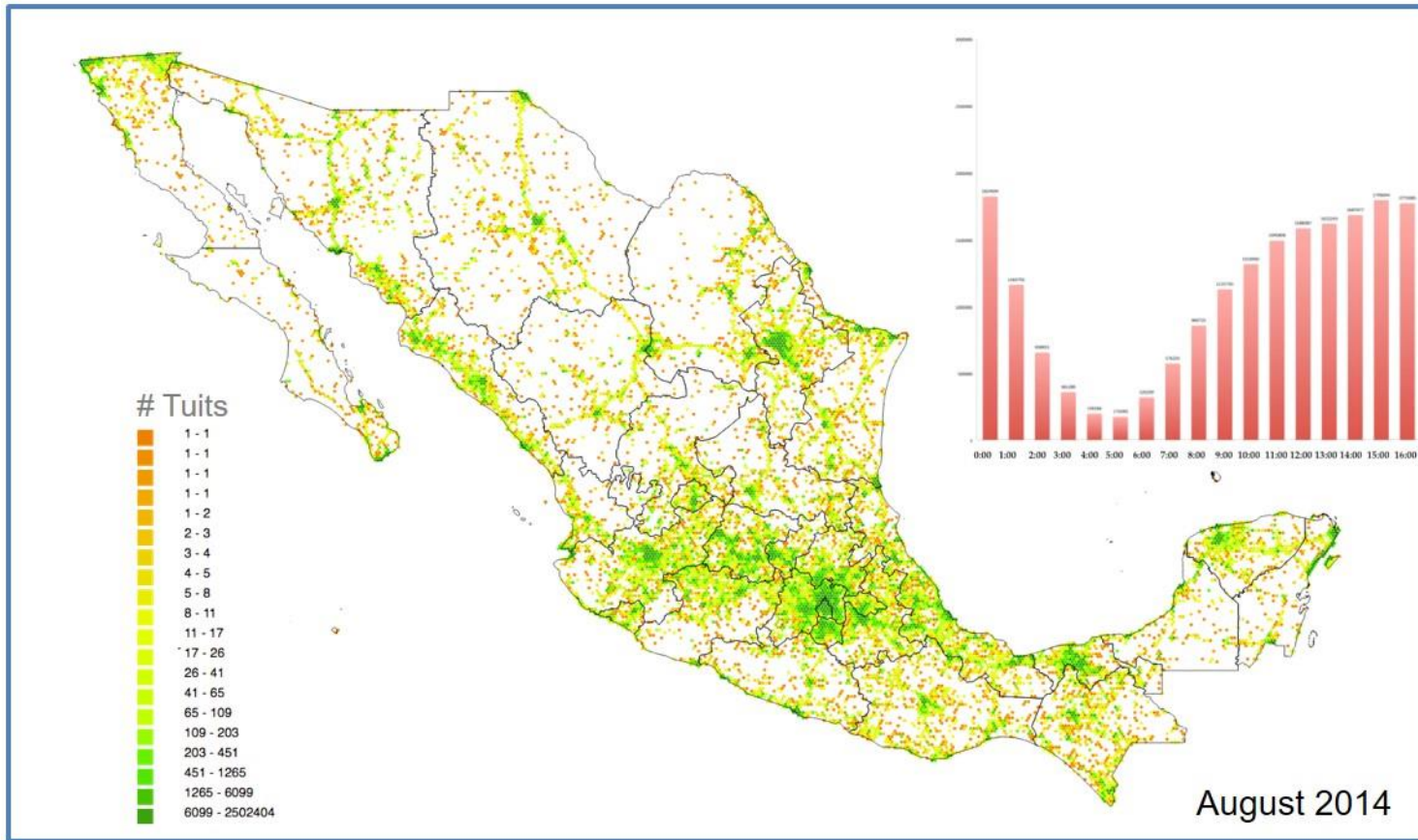
**Mexican (red) and US(blue) tweets**



**Mexican tweets**



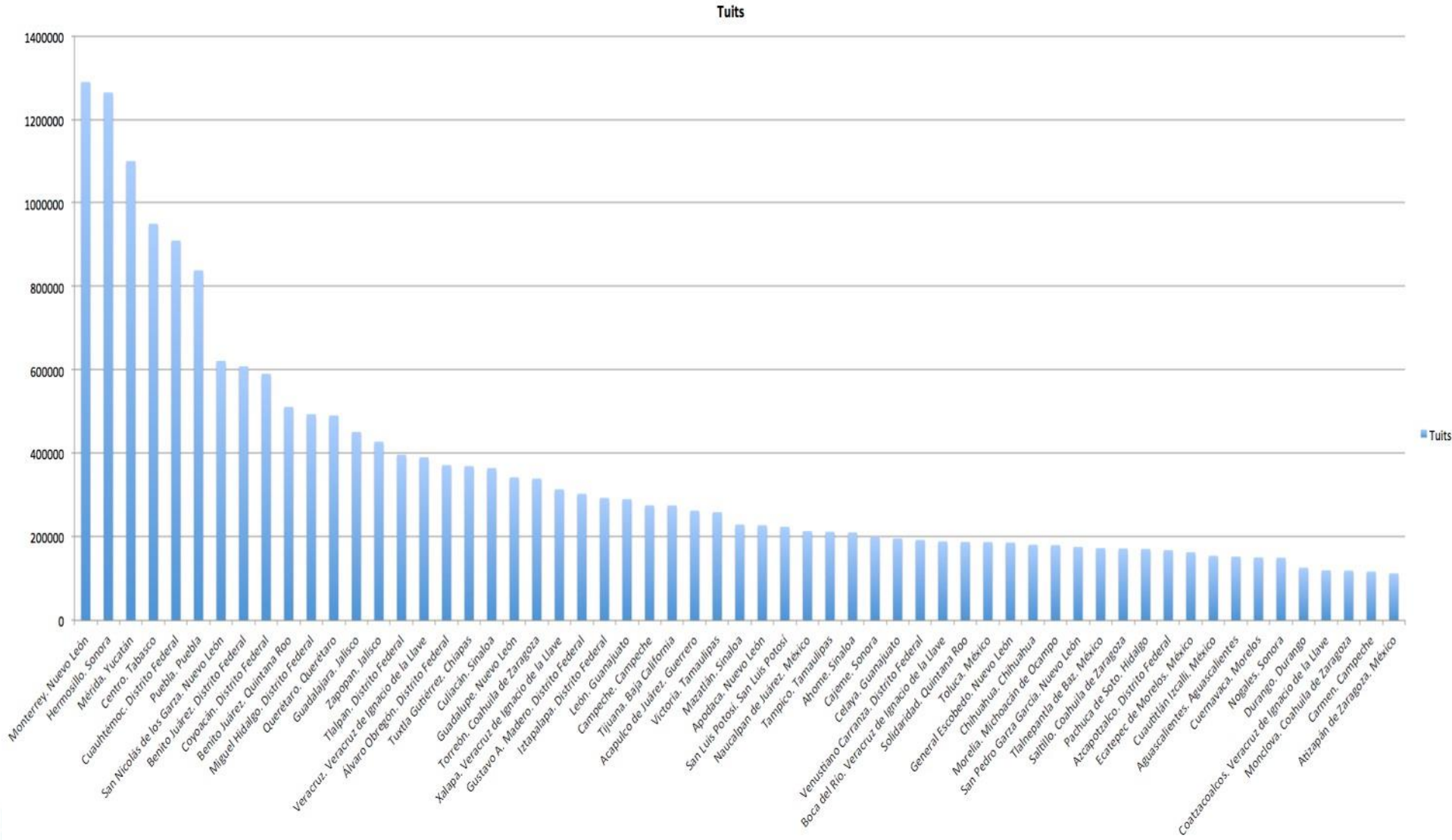
# Frequency of Tweet Generation in the country



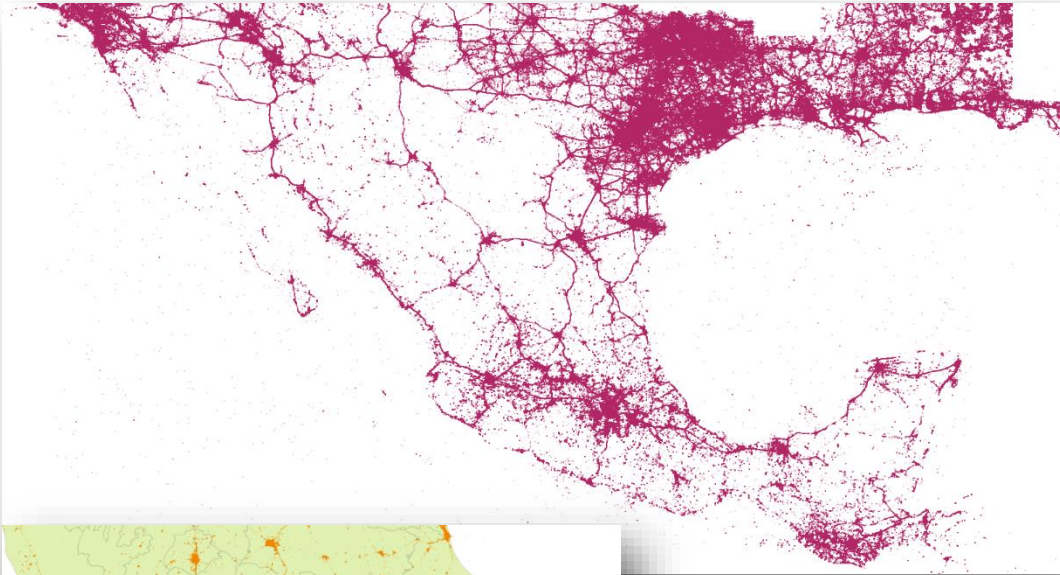
882,007 Twitter Users  
43'079,312 Geo-Referenced Tweets



# Geo-referenced Tweets in the Municipalities

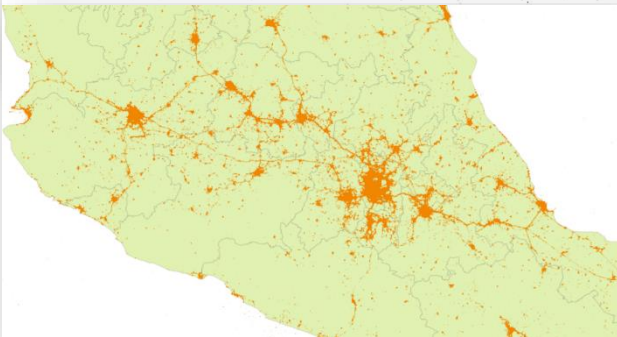


# Use of the Mexican Roads Network



121 millions of geo-referenced tweets

January 2015



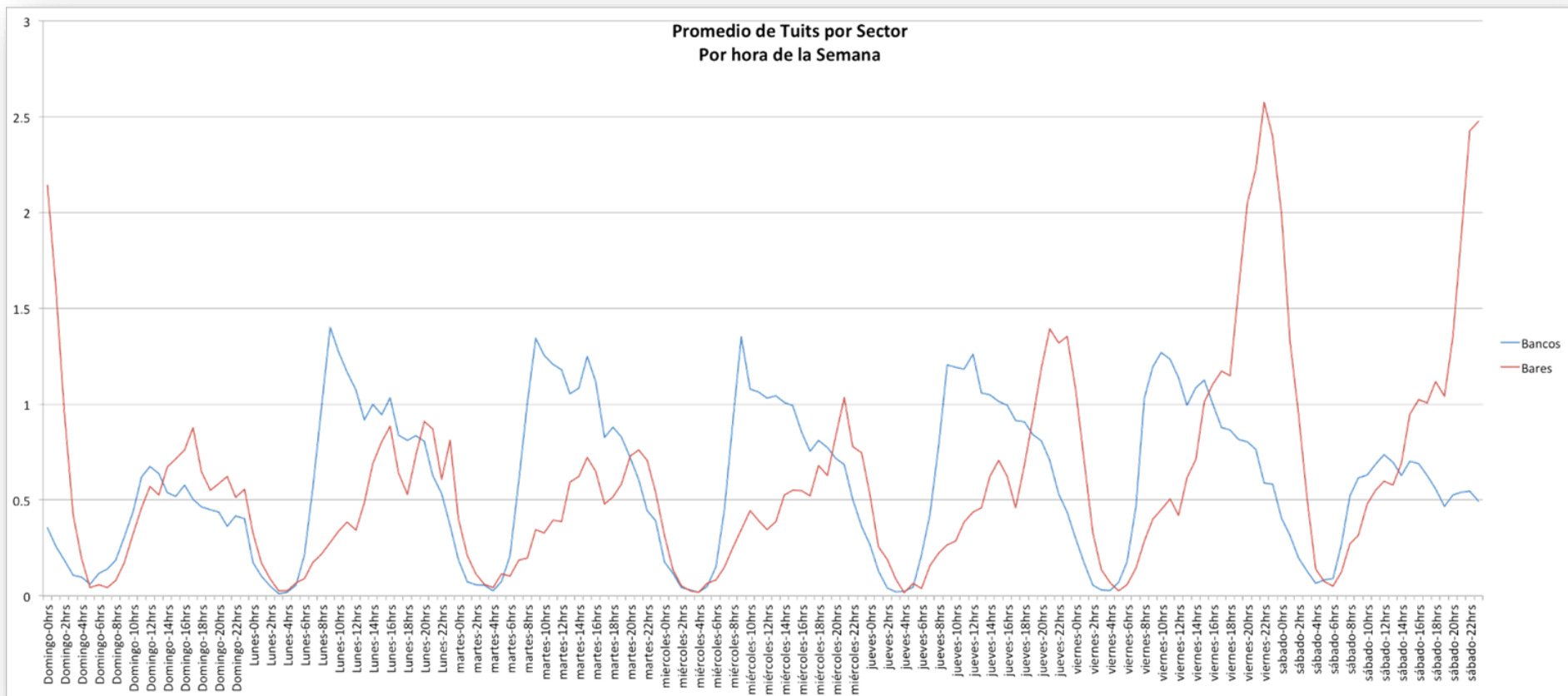
70 millions of geo-referenced tweets

October 2014



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# Average Tweets on Banks and Bars in a Week



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# Next studies on mobility...

- “Feria Nacional de San Marcos” visitants
- Internal tourism in all the country
- Mobility in our Borderlines
- Urban-Rural Systems





# SUBJECTIVE WELLNESS



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# Project Objective

**Generation of experimental indicators, new or complementary to traditional methods using data science technologies for the extraction, storage, processing and visualization of big data.**



# Expected benefits

- **New indicators obtained from Big Data Sources**
- **Correlation of results with traditional methodologies information**
- **Scientific production**



# Process

## Inception

- **Supervised Learning Method**
  - Humans put qualifications on a training set



- The system uses similarities to qualify other tweets

## Knowledge gathering...



<http://cienciadedatos.inegi.org.mx/pioanalysis>

INEGI

Desarrolló

UNIVERSIDAD TECMILENIO.  
Innovación con propósito de vida.

**PIOANALISIS** (5000+ students)

Bienestar Subjetivo y las redes sociales como fuente de datos



# Knowledge

## Set of tagged tweets from PioAnálisis (Tweet Analysis)

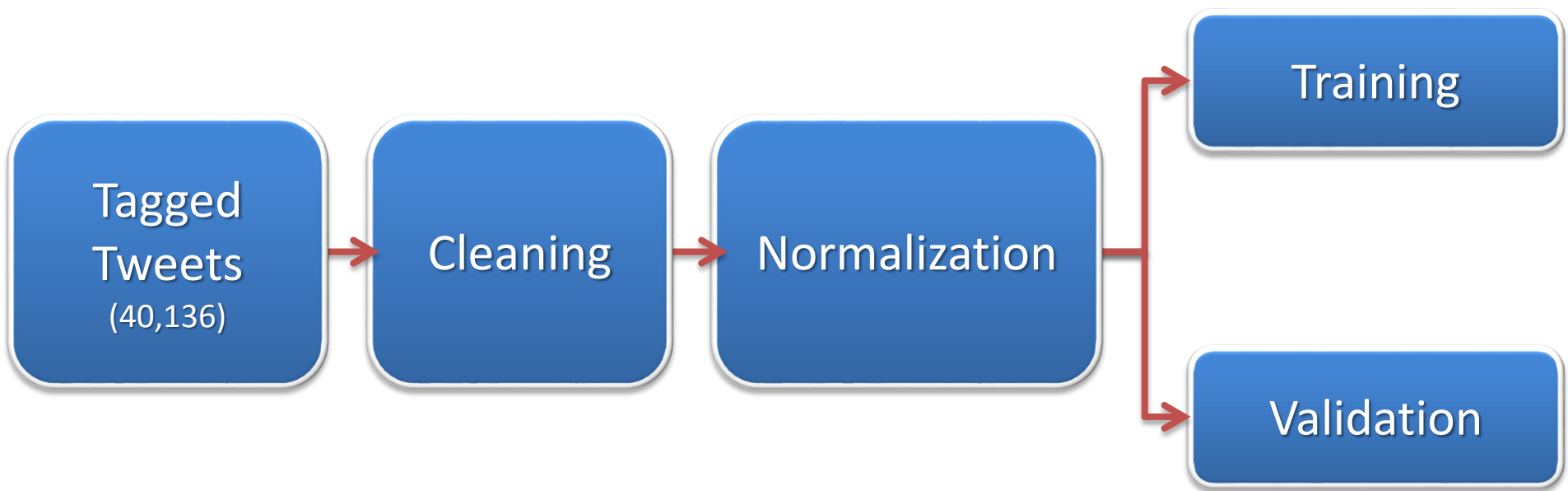
- 5000+ volunteers
- 374 000 tagged tweets
- 40 000 different tagged tweets (each tweet have been revised 9 times on average)



# Automatized Analysis and Qualification

With the manually tagged set of tweets we built a training set to teach the system to recognize and use similarities to qualify other tweets





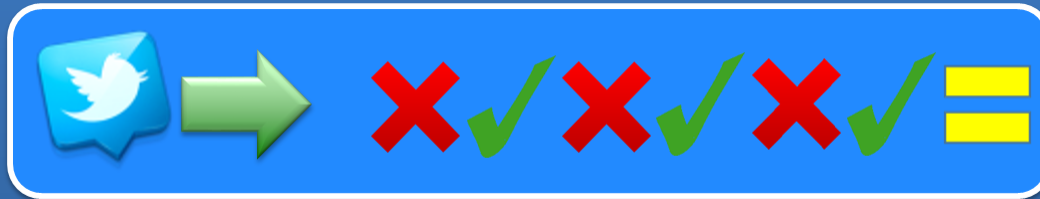


# Cleaning process

## Cleaning

Contradictions and repetitions

entropy



# Information Pre-Processing (normalization)

## Word correction

(dictionary, statistics and heuristics)

rojooooooooo jajaja

rojo

## Lematizing

(FreeLing)

rompí llegó subió

romper llegar subir

## Coding

(diacritic symbols)

caña perdón amigo

caña perdon amigo

## Polarity of Emoticons

(tag of polarity)



+1

0

-1

## Q-Grams

(q=4)

romper perder amigos

romp perd amig

## Filtering

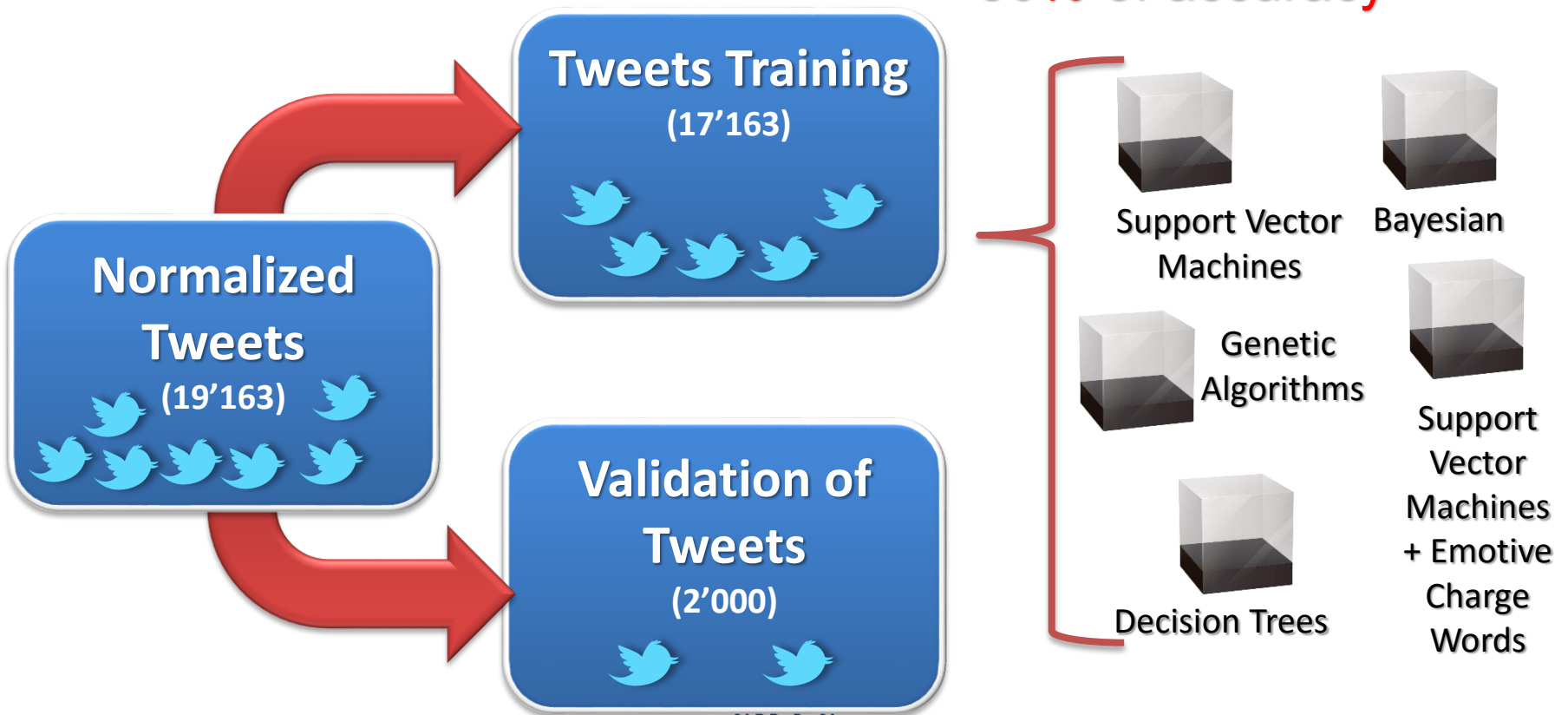
(Suppressing of stop words)

adverbs  
interjections  
verbs  
adjectives  
nouns  
hashtags

# Training with algorithms from the state of the art



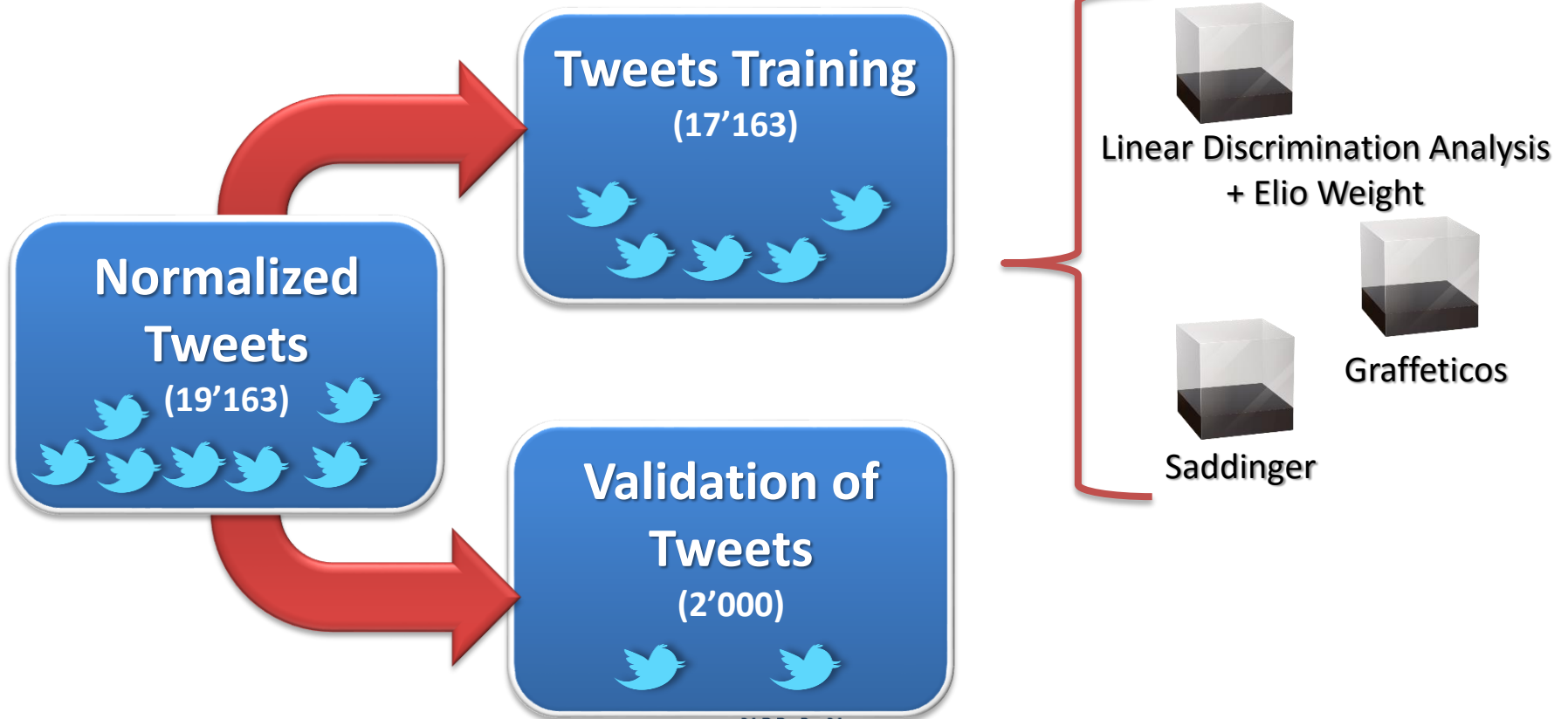
≈ 50% of accuracy



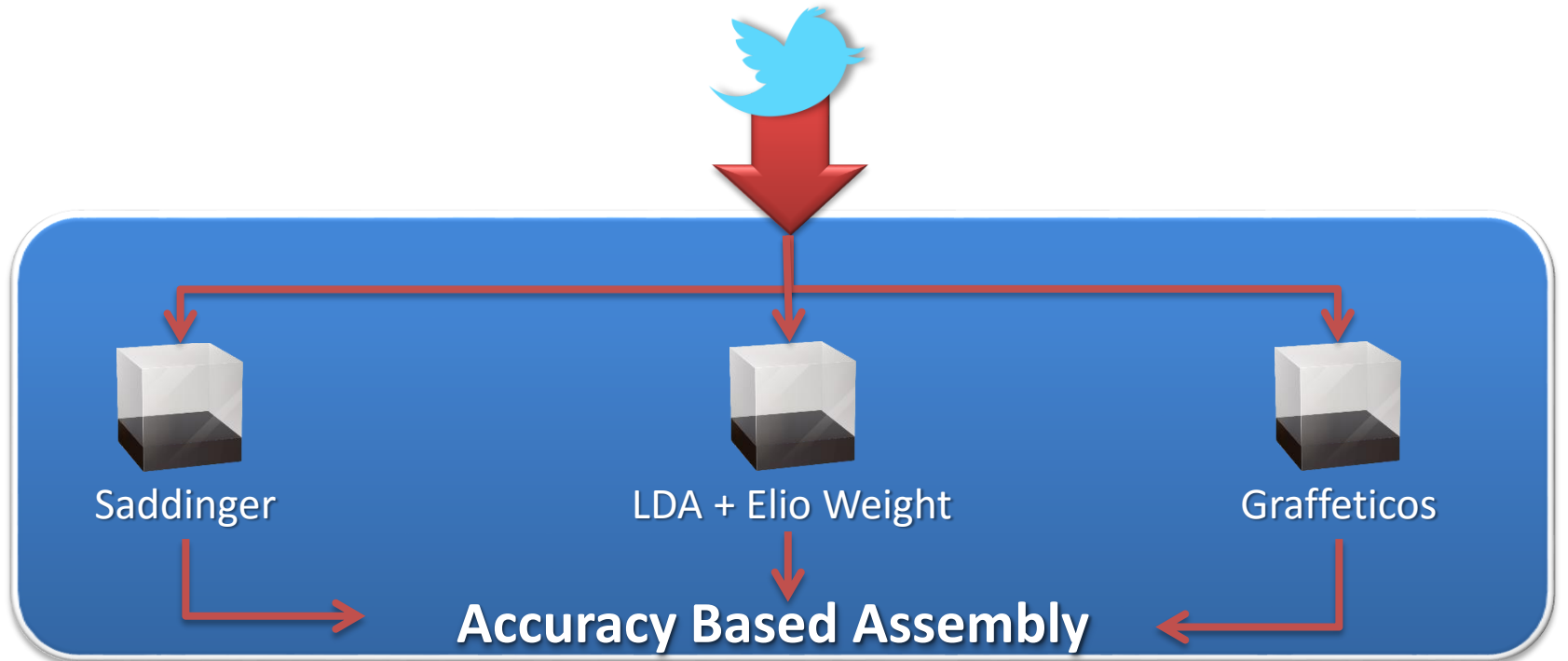
# Development of new Classification Algorithms



≈ 70% of accuracy



# Classification Algorithms Assembly to Improve Accuracy



Positive

Negative

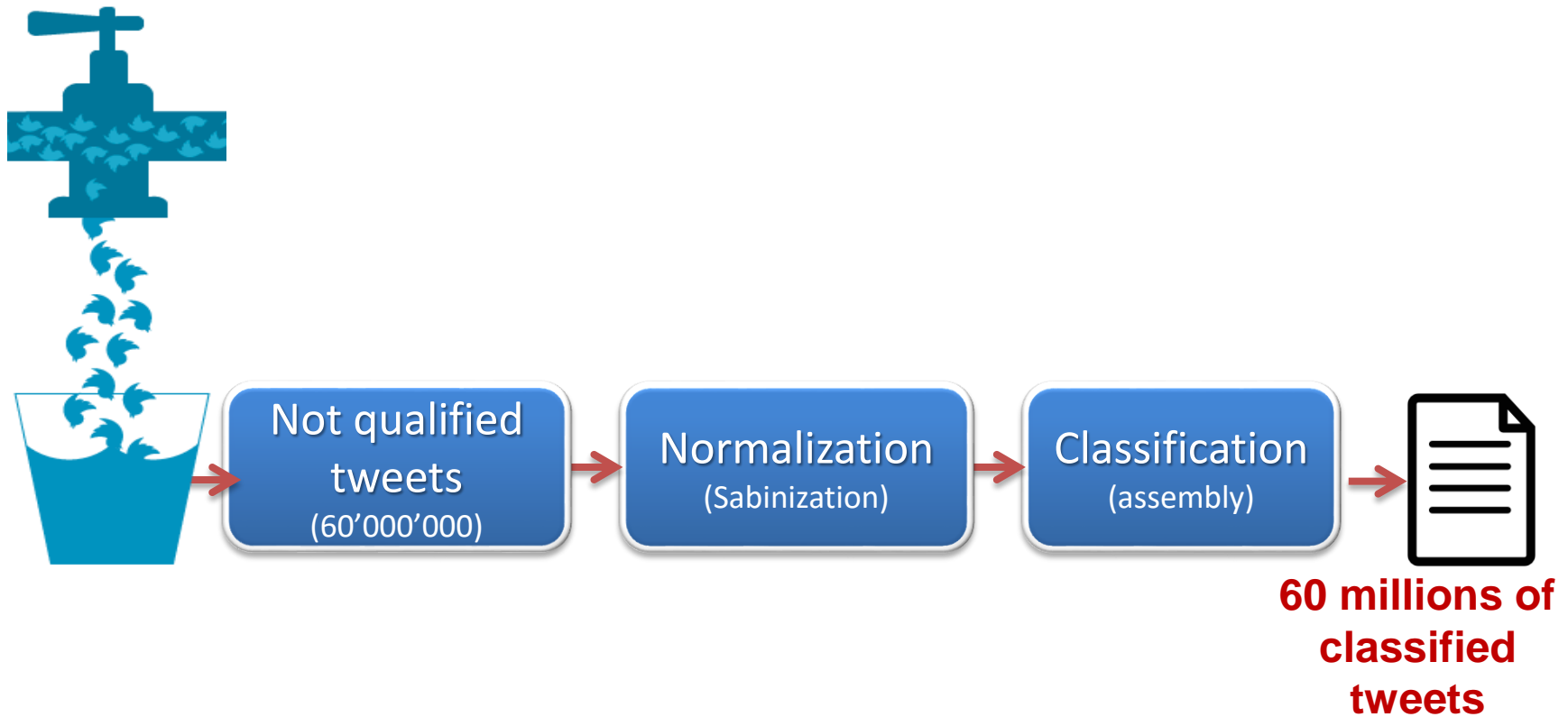


≈ 80% of accuracy



# Automatic classification of other tweets

## Classification of 60 millions of tweets with the Assembly





On line demo: [www.ingeotec.mx](http://www.ingeotec.mx)

The screenshot displays a web browser window with the URL [www.ingeotec.mx/CHA/](http://www.ingeotec.mx/CHA/). The main content is a map of Mexico and Central America, showing state and national boundaries, major cities, and geographical features. A sidebar on the left contains the logos for INFOTEC, CENTRO GEO, and INEGI. Below the logos, there is a section titled "Clasificador de sentimientos:" with a dropdown menu set to "Ensamble" and a "Reiniciar Mapa" button. At the bottom of the map area, a white rounded rectangle contains the text "Dibuje un cuadro en el mapa para iniciar la descarga de tweets...". The Google logo is visible in the bottom left corner of the map area.



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# Visualization

<http://cienciadedatos.inegi.org.mx/animotuitero/index.html>

## Estado de ánimo de los tuiteros en México

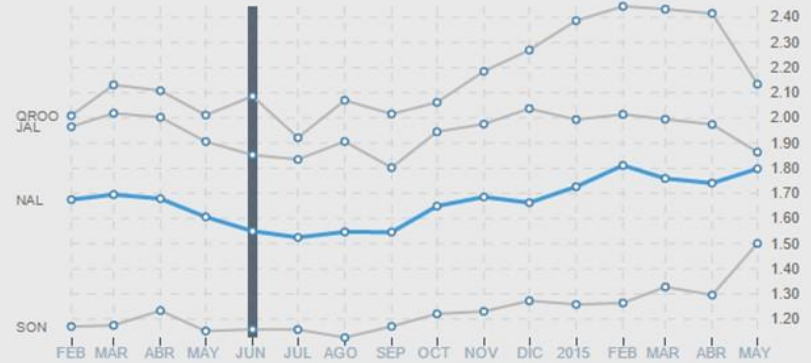


2.4



1.1

Seleccionar todo



Índice = (❤️/👎)



En colaboración con:



# High Level Meeting INEGI, INFOTEC, Centro-GEO, CIDE and CIMAT



# Results: Collaboration INEGI, INFOTEC, Centro-GEO, CIDE y CIMAT

- Collaboration lines:
  - Common research
  - Seminars
  - Internships
  - academic programs
  - Spaces exchange
  - Micro data access



# NEW PROJECTS



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

# Social Networks monitoring for INEGI

**Goal:** Publication, following and monitoring of social networks

- Evaluation of new tools: Semantic Web Builder developed by INFOTEC
- Living Lab workshop: Dissemination staff
- Implementing on internal environment





# Mental Health of Teenagers (Data2X)

**Objective:** “Generation of information about mental health on teenager women in Mexico from Tweeter messages”

- INEGI-Data2X agreement (Data2X is a ONG supported by UN)



# Mental Health of Teenagers (Data2X)



# Statistics on Security and Safety

Research on the possibility to use the tweets database to get information about:

- Collection of data in urban areas
- Information about natural disasters





## Conociendo México

01 800 111 46 34

[www.inegi.org.mx](http://www.inegi.org.mx)

[atencion.usuarios@inegi.org.mx](mailto:atencion.usuarios@inegi.org.mx)

# Thank you!!



[@inegi\\_informa](https://twitter.com/inegi_informa)



[INEGI Informa](https://www.facebook.com/INEGIInforma)



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

[Juan.Munoz@inegi.org.mx](mailto:Juan.Munoz@inegi.org.mx)

