

Sandbox Project On Social Media Data for Sentiment Analysis/Mobility

3rd International Conference on Big Data for Official Statistics
30 Aug – 1 Sep 2016



INSTITUTO NACIONAL
DE ESTADÍSTICA Y GEOGRAFÍA

 **BigData** UN Global Working Group

Objective

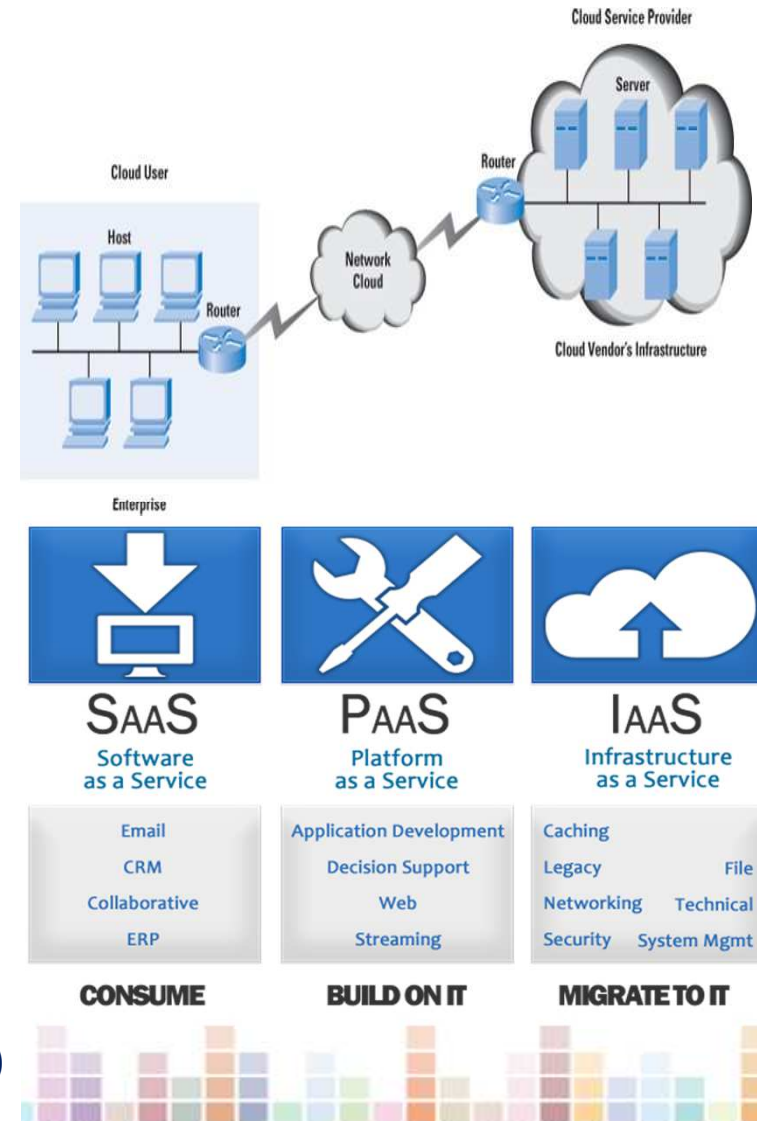
To share the experience of INEGI in the use of Twitter as a Big Data source and his collaboration in the Sandbox project



INSTITUTO NACIONAL
DE ESTADÍSTICA Y GEOGRAFÍA

Cloud Computing

Cloud computing is a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction.



US National Institute of Standards (NIST)

Sandbox, since 2014

- Braced by the Commission of European Statisticians (CES)
- Promoted by High Level Group for the Modernization of Official Statistics (HLG-MOS)
- Coordinated by the Statistical Division of the United Nations Economic Commission for Europe (UNECE)
- Implemented in the Irish Centre for High-End Computing



Sandbox Initial Aims

- Test feasibility of remote access and processing: Could this approach be used in practice?
- Test whether existing statistical standards / models / methods can be applied to Big Data
- Determine which Big Data software tools are most useful for statistical organisations
- Learn about the potential uses, advantages and disadvantages of Big Data – “learning by doing”.
- Build an international collaboration community on the technical aspects of using Big Data

modern
stats



What can we expect from the Sandbox?

- A state-of-the-art shared computing environment, or in other words: a high-end collaborative ICT environment to make experiments
- Hardware:
 - 4 Data/compute nodes, each with 2x10 core Intel Xeon CPUs, 128 GB RAM, 4x4TB disks, 56 Gbit InfiniBand network
 - 2 Service/login nodes , each with 2x10 core Intel Xeon CPUs, 128 GB RAM, 4x4TB disks, 10 Gbit connection to Internet
- Software:
 - Hortonworks Data Platform (Hadoop, Spark, Hive, Pig), R-Studio, RHadoop, ElasticSearch



INSTITUTO NACIONAL
DE ESTADÍSTICA Y GEOGRAFÍA

Types of Big Data Sources

- **Meters and smart sensors.** Traffic cameras, GPS devices, power consume meters, IoT, smartwatches, smartphones, etc.
- **Social interactions.** Conversations and publications on social networks like Twitter, Facebook, FourSquare, etc.
- **Business transactions.** Credit cards movements, scanned data, cell phone records, etc.
- **Electronic files.** Documents which are available in electronic formats such as PDF files, websites, videos, audio, images, photos, etc.
- **Broadcast media.** Digital video and audio streamed on real-time



INSTITUTO NACIONAL
DE ESTADÍSTICA Y GEOGRAFÍA

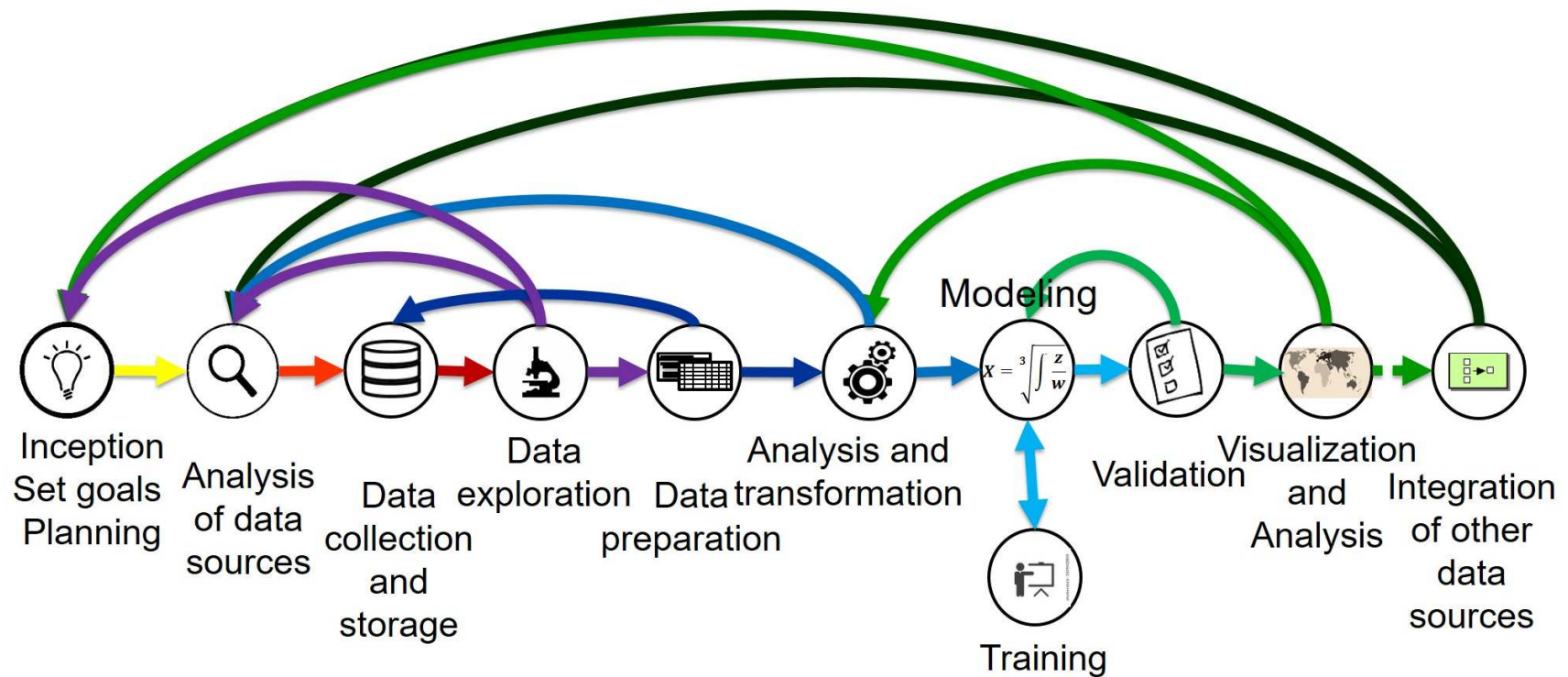
Types of Big Data Sources

- **Meters and smart sensors.** Traffic cameras, GPS devices, power consume meters, IoT, smartwatches, smartphones, etc.
- **Social interactions.** Conversations and publications on social networks like Twitter, Facebook, FourSquare, etc.
- **Business transactions.** Credit cards movements, 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 streamed on real-time

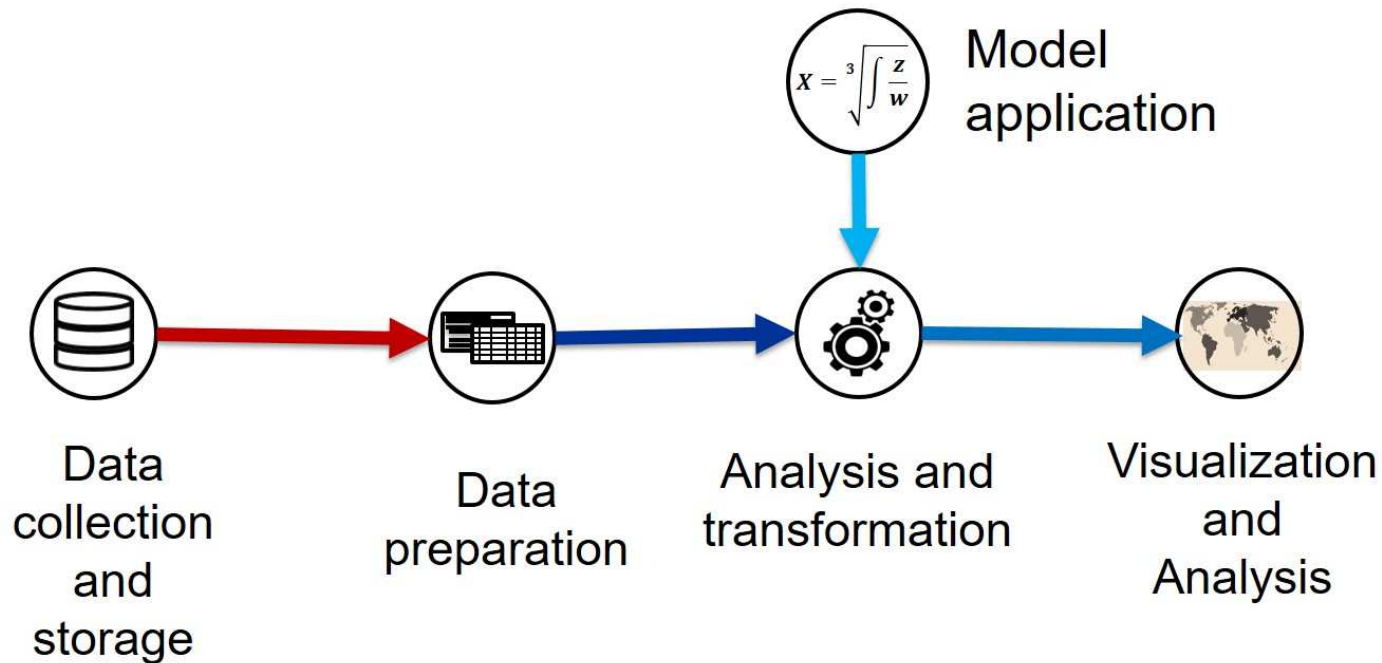


INSTITUTO NACIONAL
DE ESTADÍSTICA Y GEOGRAFÍA

Process Followed by INEGI (until now)



Production process



INSTITUTO NACIONAL
DE ESTADÍSTICA Y GEOGRAFÍA

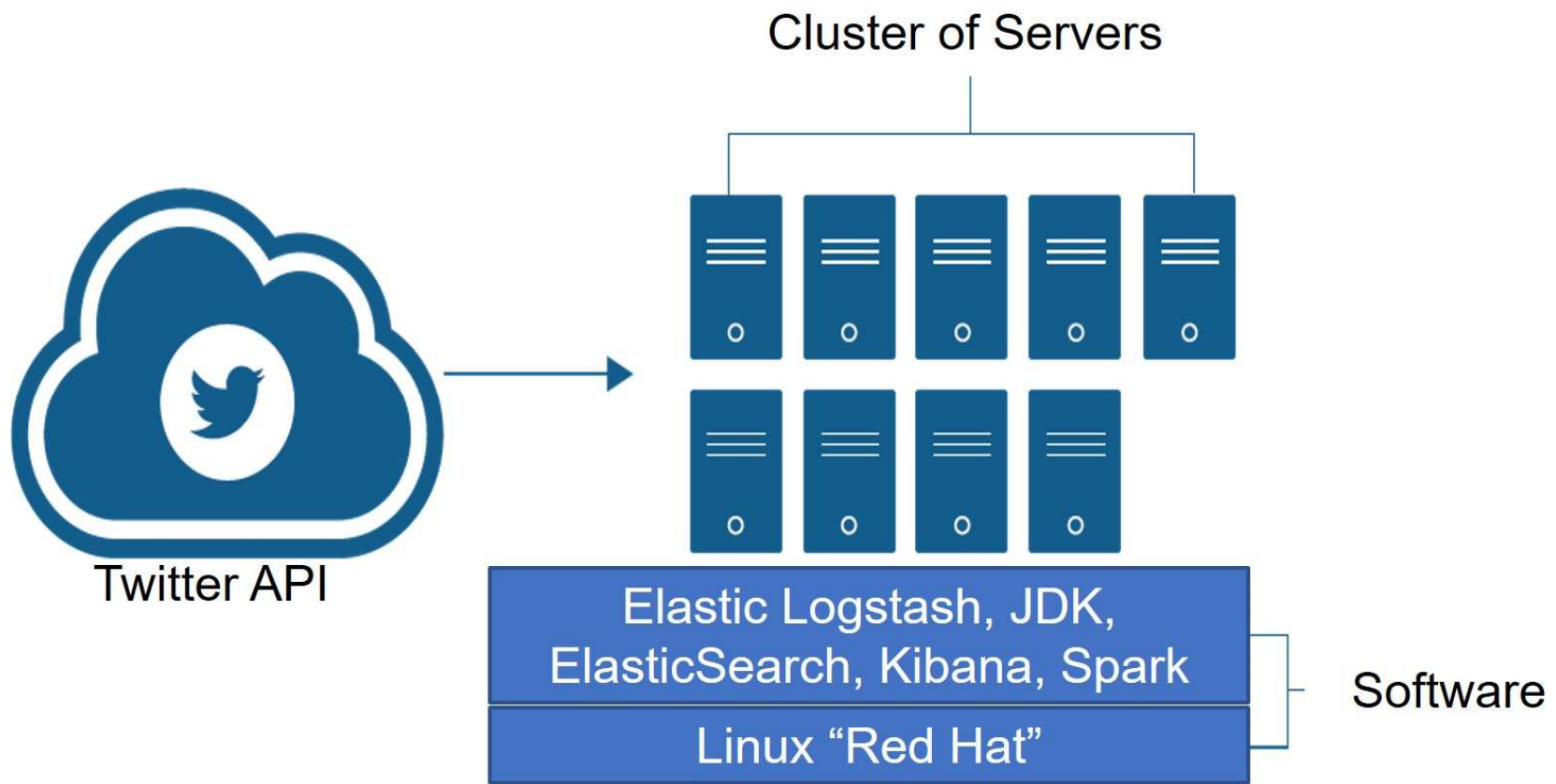
Study Case

- **Initial objective of INEGI's Big Data Project:** To generate experimental indicators using Big Data techniques with social media data, to complement statistical information obtained from traditional methods and sources.
- **Initial Goal:** To obtain indicators of subjective wellbeing from social media data sources.

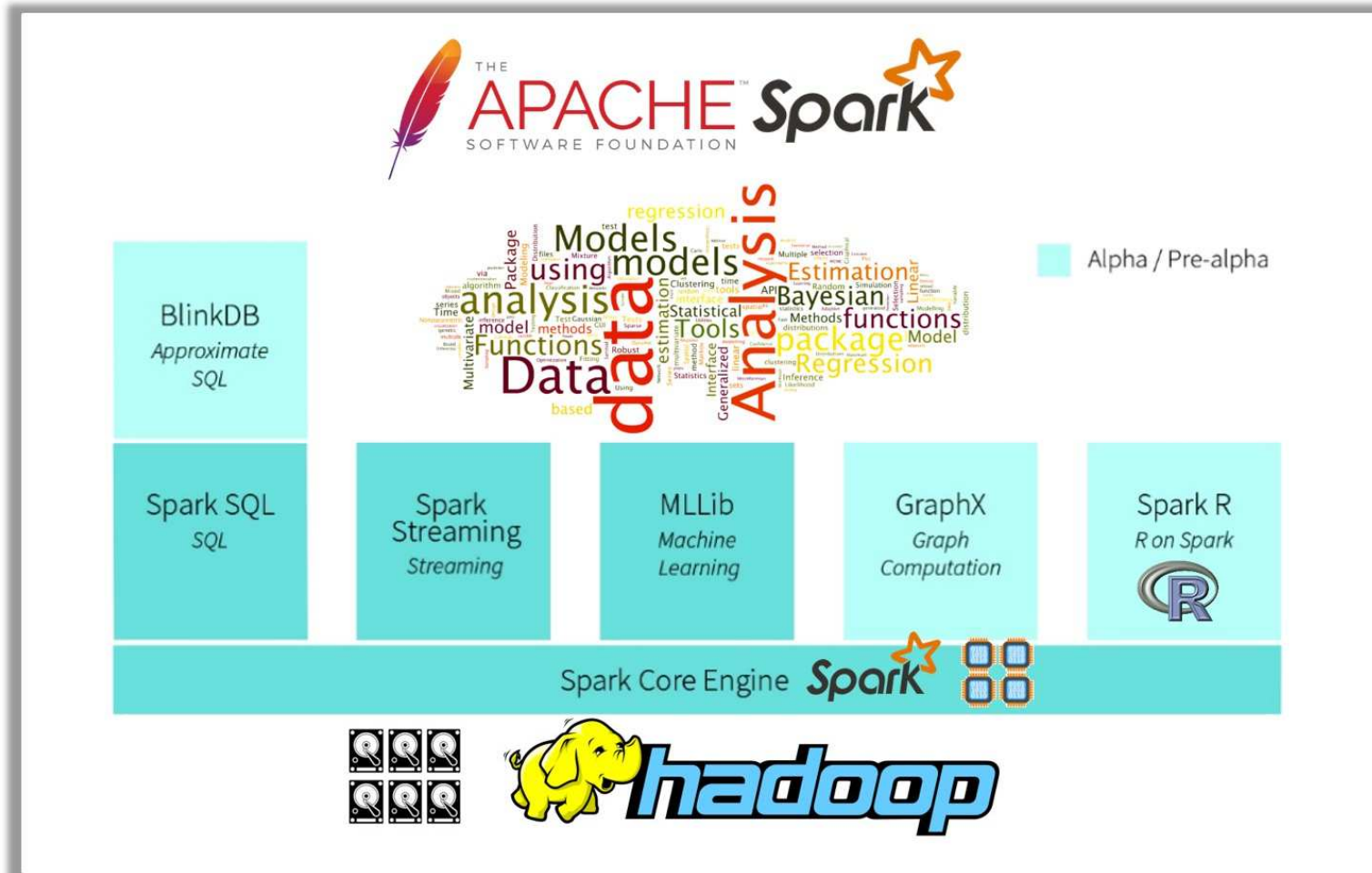
Why did we choose Twitter as a data source?

- It's a widely adopted social network where you can find content written by common people
- Tweets are public, so we can use them without concerns about privacy
- There is a free API which allows to get 1% of the tweets that are being produced on real time (<https://dev.twitter.com/rest/public>)

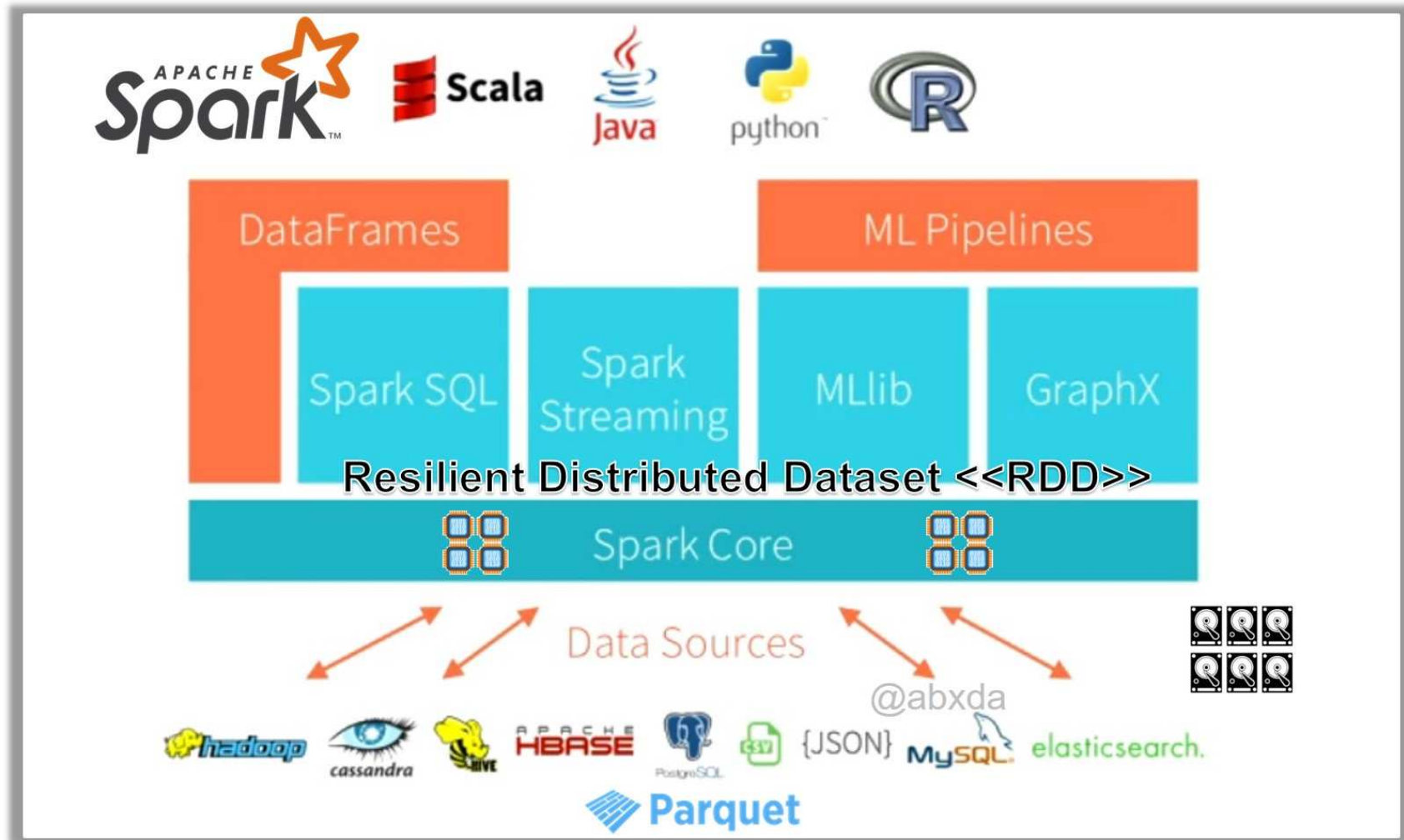
Collection infrastructure



Software Stack (2013)

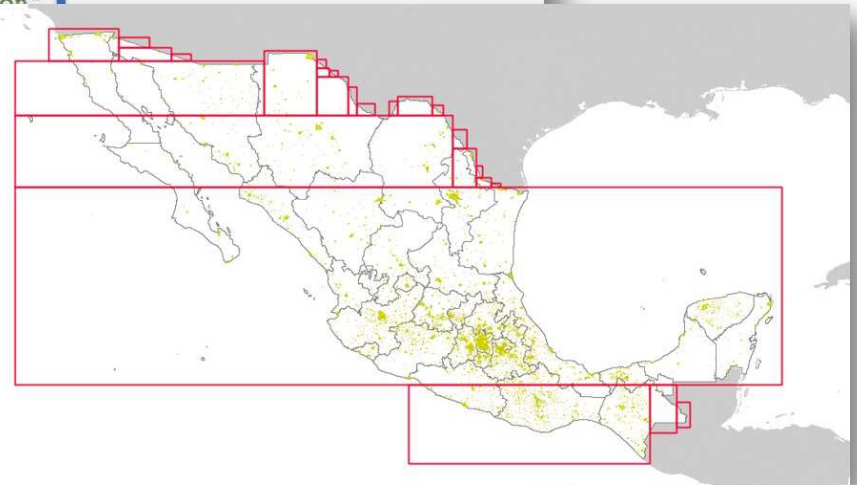


Software Stack (2016)



Tweet Structure

```
{
  "text": "Todo listo para la presentación de #BigData en el @FSLmx",
  "created_at": "2014-02-21T17:00:50.000Z",
  "truncated": false,
  "mention": [],
  "retweet_count": 0,
  "hashtag": [],
  "location": {
    "lat": 19.39617897,
    "lon": -99.22636055
  },
  "place": {
    "id": "3ad512d283f67a11",
    "name": "Aguascalientes",
    "type": "city",
    "full_name": "Aguascalientes, Aguascalientes",
    "street_address": null,
    "country": "México",
    "country_code": "MX",
    "url": "https://api.twitter.com/1.1/geo/id/3ad512d283f67a11/reverse-geo"
  },
  "link": [
    {
      "url": "http://t.co/AUNXLVSImQ",
      "display_url": "4sq.com/lplLyUL",
      "expand_url": "http://4sq.com/lplLyUL",
      "start": 28,
      "end": 50
    }
  ],
  "user": {
    "id": 205760874,
    "name": "Abel Coronado",
    "screen_name": "abxda",
    "location": "",
    "description": "Filósofo, Desarrollador de Software, M.C. en I",
    "profile_image_url": "http://pbs.twimg.com/profile_images/...",
    "profile_image_url_https": "https://pbs.twimg.com/profile_..."
  }
}
```



We found that

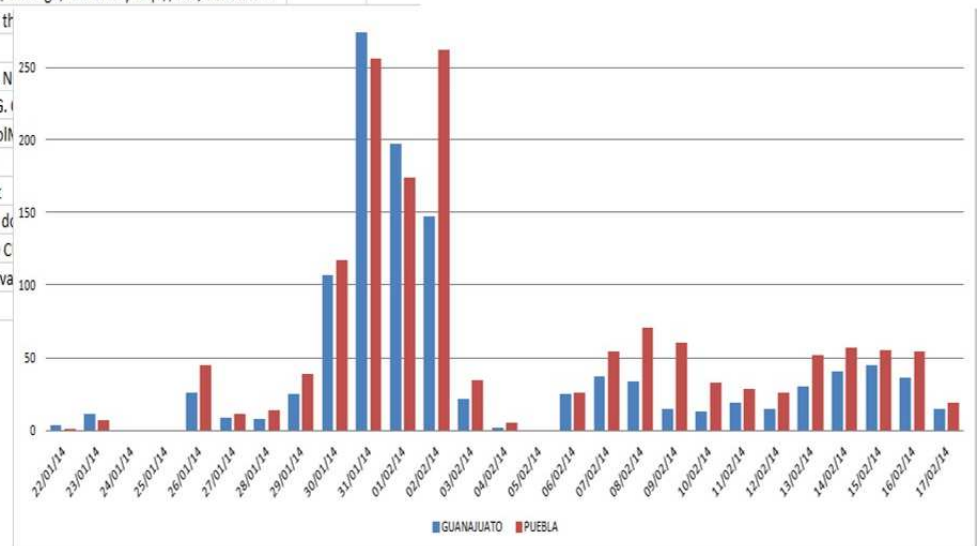
- The JSON structure is easy to process
- The content is text that we can examine to make the sentiment analysis
- Geographical coordinates can be used to filter the tweets and obtain only those of interest (warning: not all the tweets are geo referenced)
- We can make mobility analysis, based in Tweets' location and time



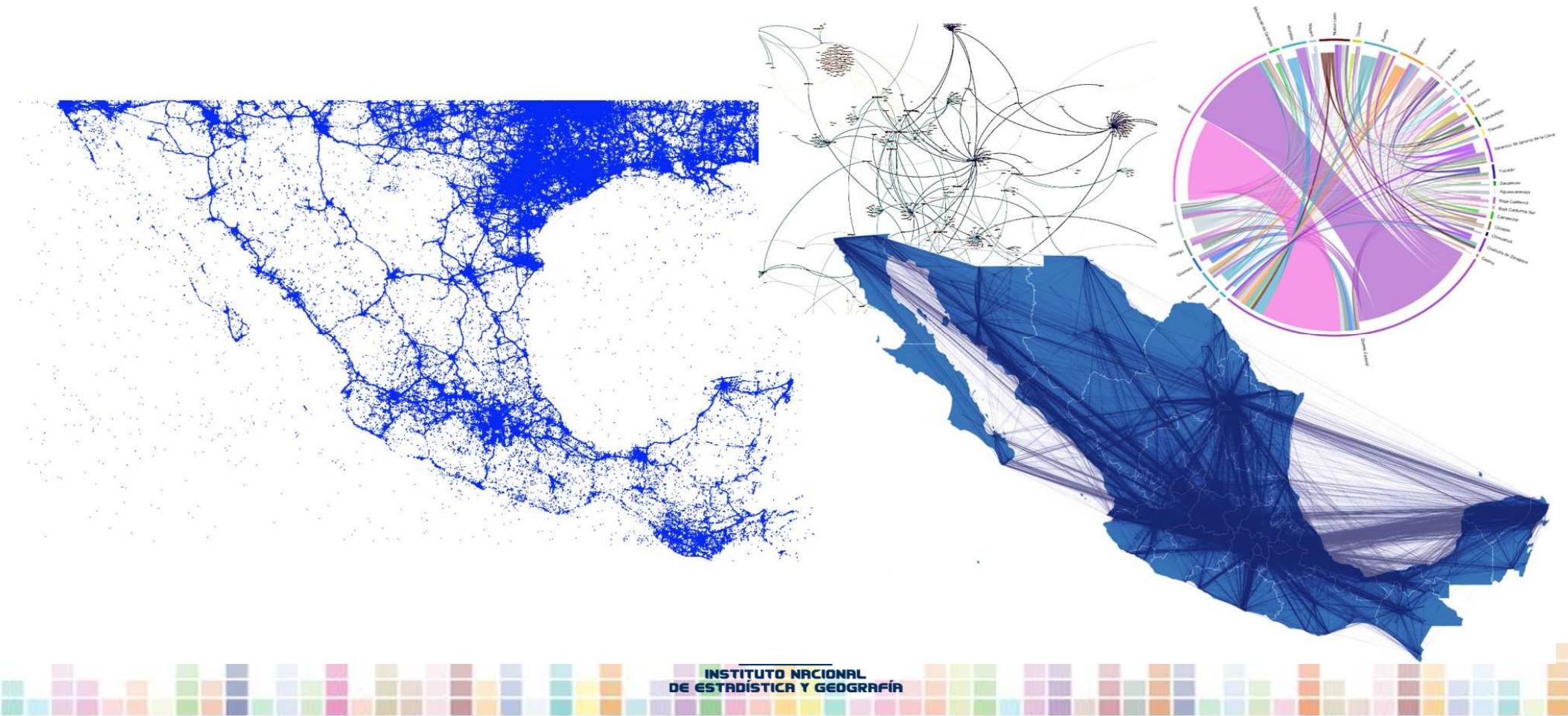
INSTITUTO NACIONAL
DE ESTADÍSTICA Y GEOGRAFÍA

Tweets preparation and analysis

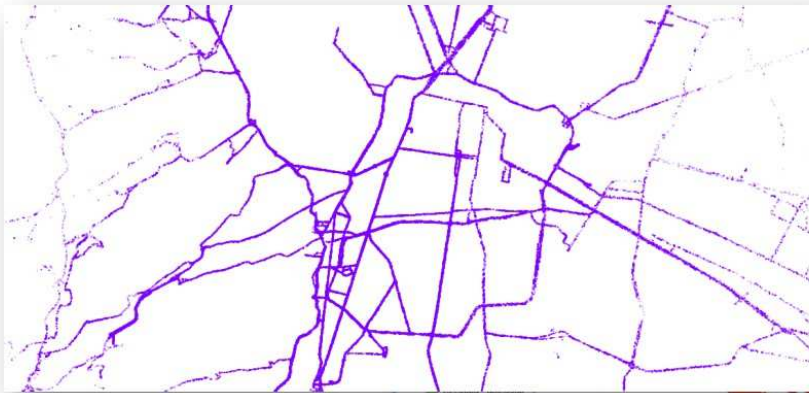
4.3691E+17	ayyekaylla	30.792072	-86.67873	US	<a href="htt
4.3691E+17	_OnlyOneJa	31.328173	-89.329953	US	<a href="htt
4.3691E+17	ww_cesar_w	19.396179	-99.226361	MX	<a href="htt
4.3691E+17	Ericikee	27.391716	-82.477106	US	<a href="htt
4.3691E+17	jennylove	29.709114	-95.185322	US	<a href="htt
4.3691E+17	fets	26.073966	-80.147174	US	<a href="htt
4.3691E+17	brettwbailey	30.109563	-91.942463	US	<a href="htt
4.3691E+17	MiOfiPolanc	19.435647	-99.185377	MX	<a href="htt
4.3691E+17	kleptolovestory			US	<a href="htt
4.3691E+17	maxoyervide	26.211199	-98.21623	US	<a href="htt
4.3691E+17	balehorabue	25.671282	-100.313277	MX	<a href="htt
4.3691E+17	YoaniitaFran	21.144447	-86.906771	MX	<a href="htt
4.3691E+17	nicole_rendi	24.02279	-104.672112	MX	<a href="htt
4.3691E+17	MikelCass	30.699126	-97.864629	US	<a href="htt
4.3691E+17	alezavaliis	15.466293	-87.996081	MX	<a href="htt
4.3691E+17	carlyy_g	30.199626	-97.773748	US	<a href="htt
4.3691E+17	ElColver			MX	web
4.3691E+17	lrobertonav	23.194018	-106.425074	MX	<a href="htt
4.3691E+17	SheFellForM	32.558248	-97.057816	US	<a href="htt
4.3691E+17	EhmJayyOrN	30.423385	-91.173801	US	<a href="htt
4.3691E+17	_Kenahday_l	32.500404	-86.471853	US	<a href="htt
4.3691E+17	HaleyMathe	30.609897	-96.333299	US	<a href="htt
4.3691E+17	enlacedelacosta			MX	web
4.3691E+17	karisalinast	25.690377	-100.276162	MX	<a href="htt



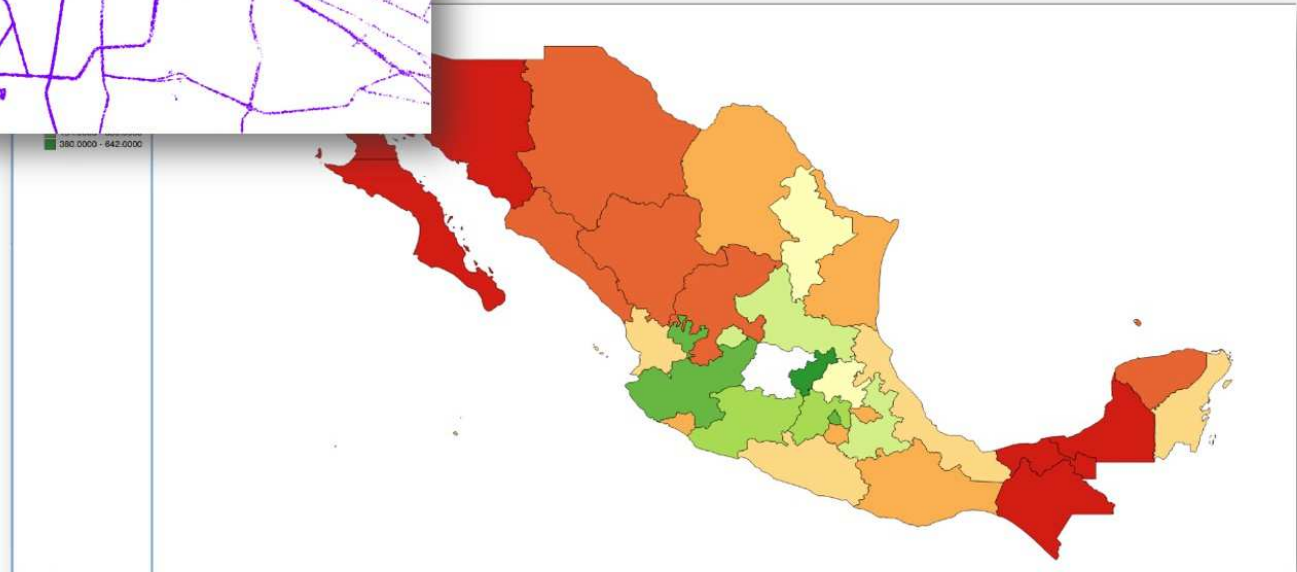
Modeling and Validation



Visualization

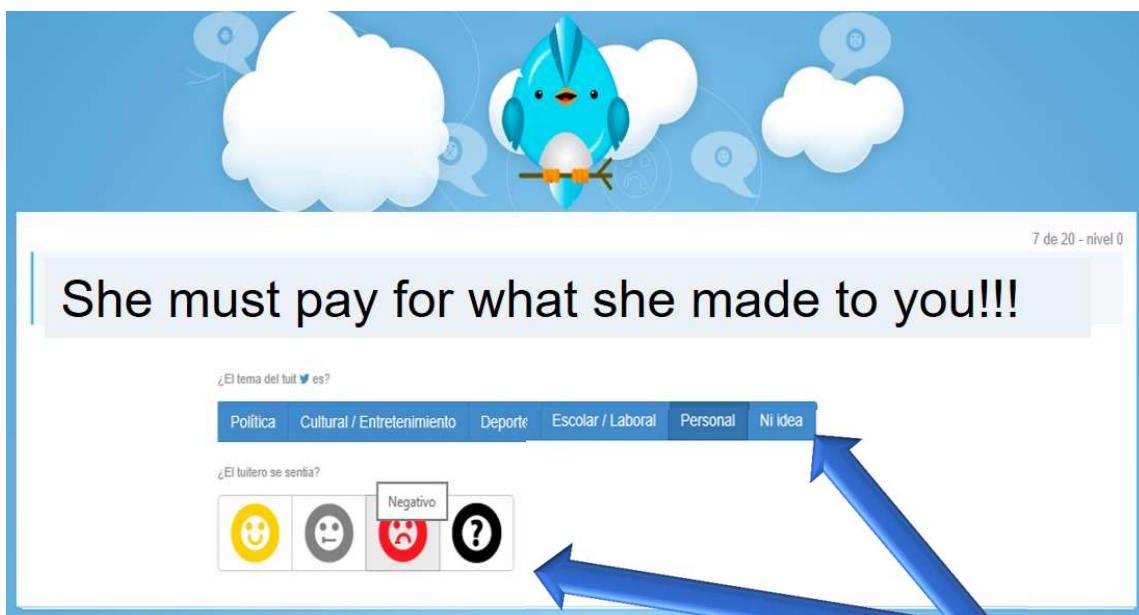


380 0000 - 642 0000



Supervised Training

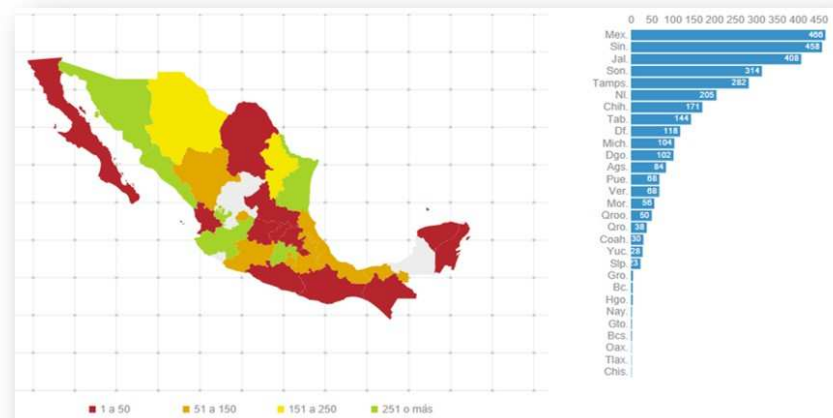
<http://cienciadedatos.inegi.org.mx/animotuitero/>



Manual Tagging



INSTITUTO NACIONAL
DE ESTADÍSTICA Y GEOGRAFÍA

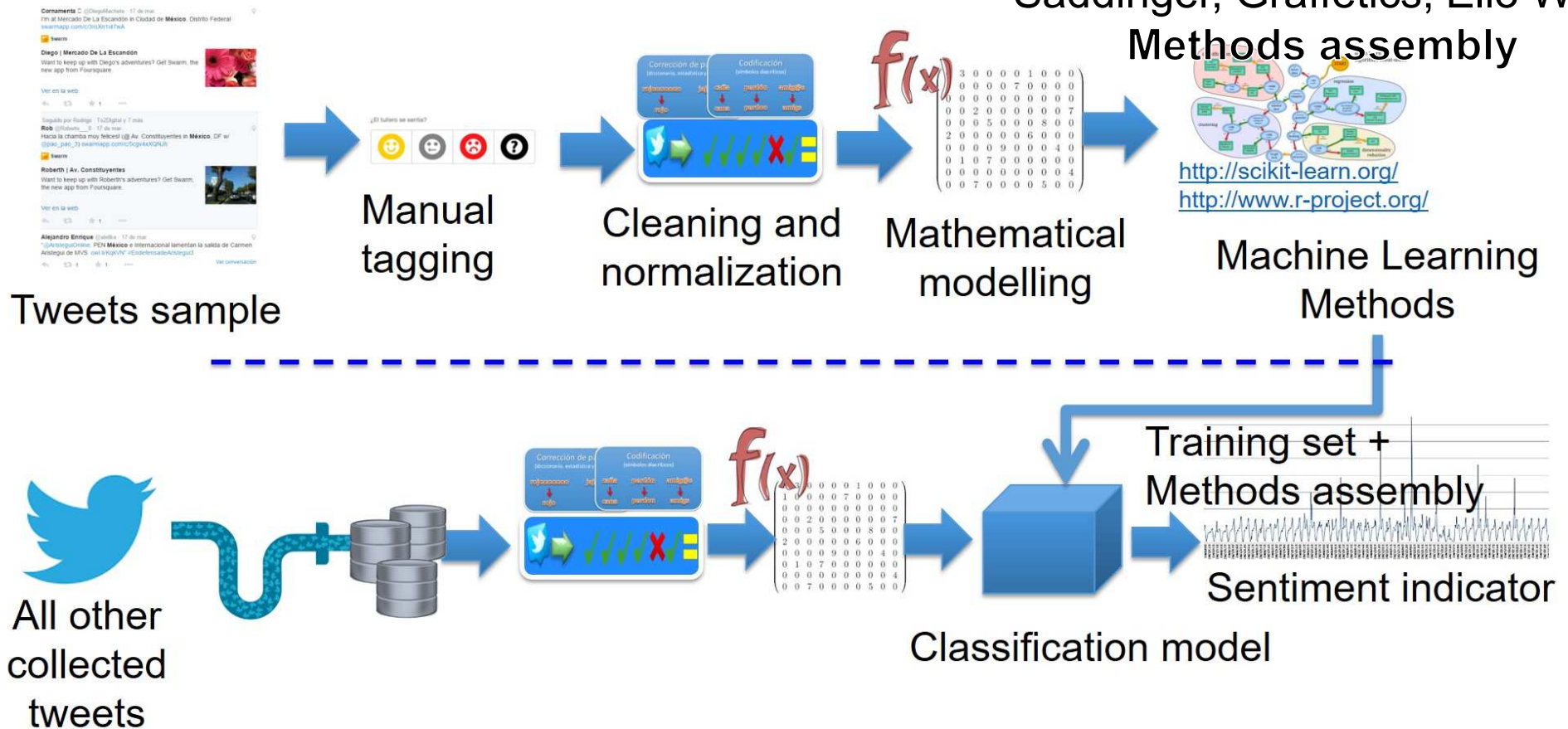


5000 people (TecMilenio students),
100 tweets tagged by each one,
each tweet was tagged nine times,
about 40,000 different tagged tweets,
interpreted accordingly to regional
idioms

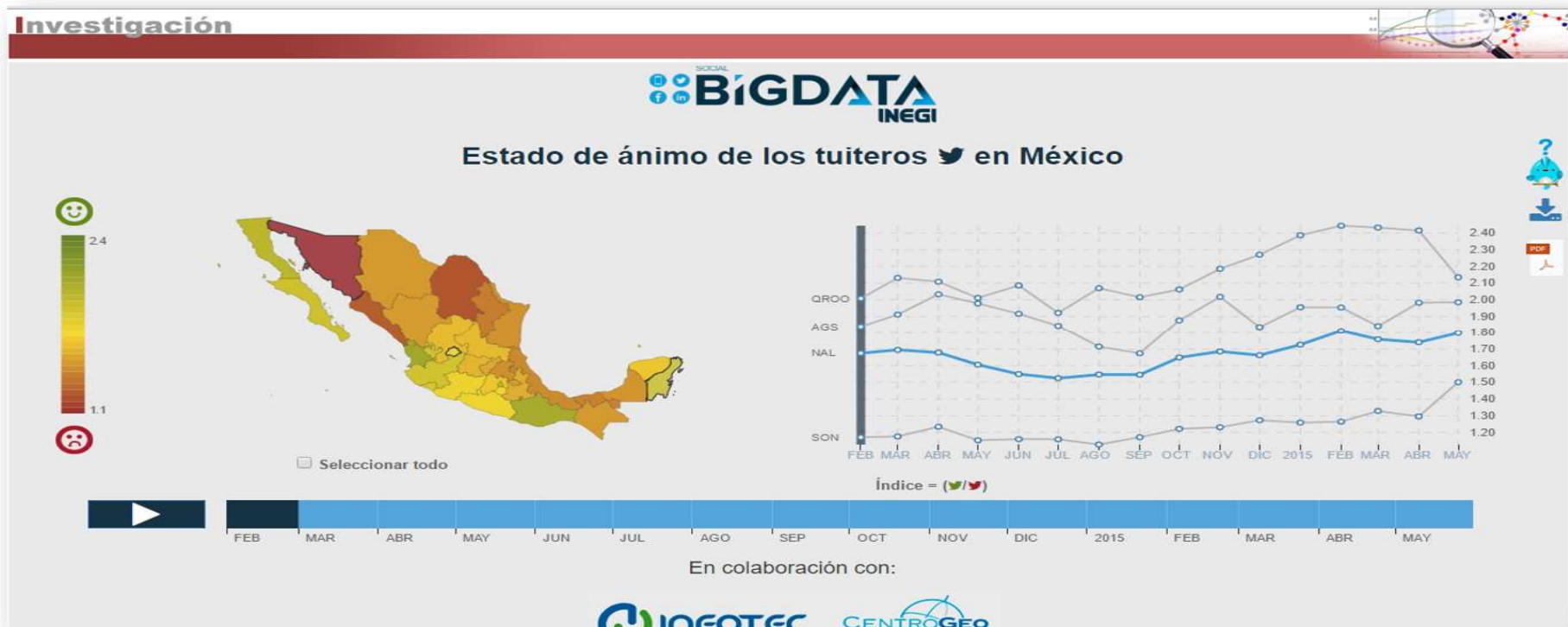
Training and modeling

- Latent Dirichlet Allocation
- Decision Trees
- Saddinger, Graffetics, Elio Weight

Methods assembly

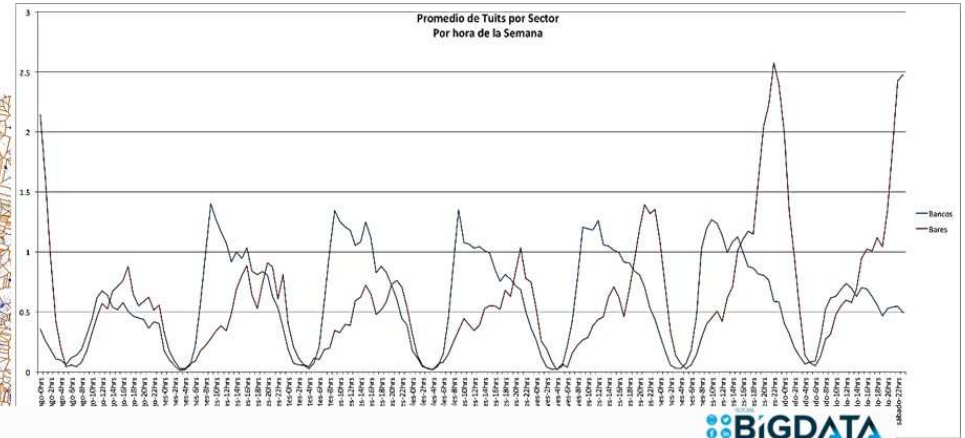
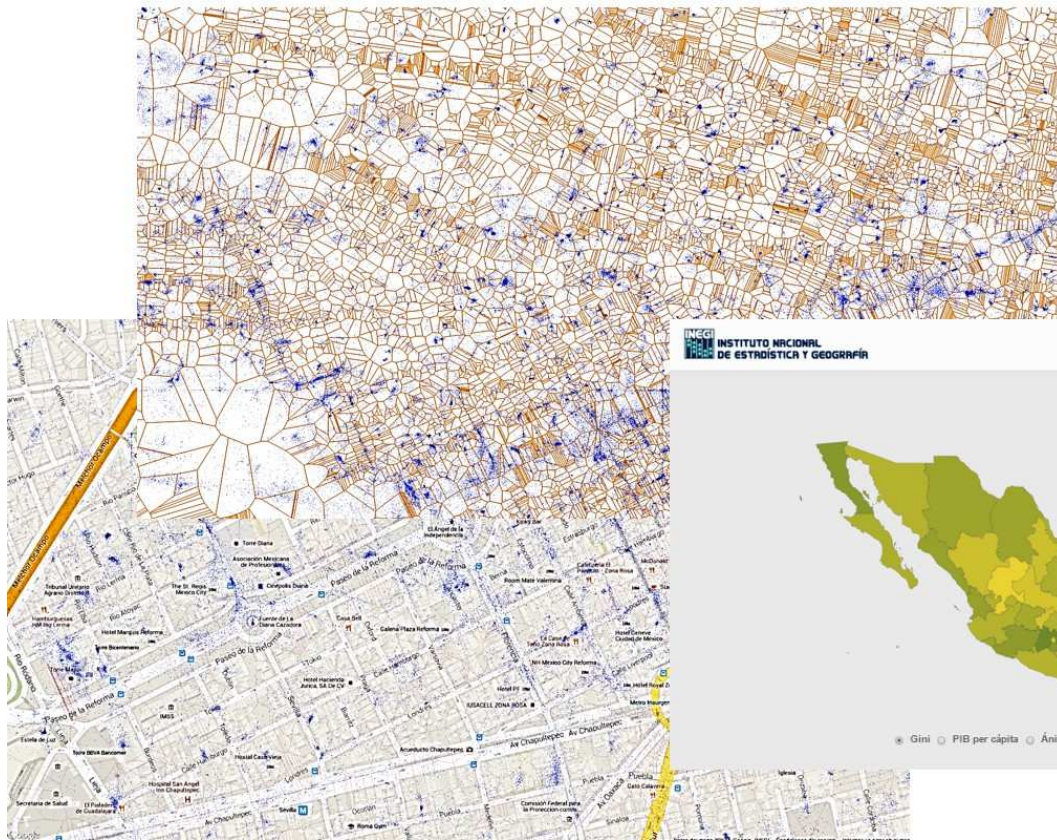


Sentiment Visualization

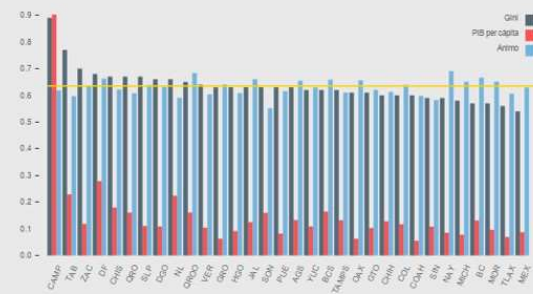


<http://www.inegi.org.mx/inegi/contenidos/investigacion/Experimentales/animotuitero/default.aspx>

Integration of other sources



Desigualdad en México 2012



Some applications

- Tourism
- Migration
- Use of roads
- Regional influence of big cities
- Mobility patterns
- Business activity patterns
- Subjective wellbeing
- Inequities impact
- Impact analysis of relevant news
- Mental health
- Misogynist/discriminatory language use
- SDG indicators?



INSTITUTO NACIONAL
DE ESTADÍSTICA Y GEOGRAFÍA

Collaboration

- International
 - UNECE
 - ICHEC
 - UNSD
 - LAMBDoop
 - University of Pennsylvania
- National
 - KioNetworks
 - Dattlas
 - TecMilenio
 - INFOTEC
 - Centro Geo
 - CIDE
 - CIMAT
 - Sector
- Internal
 - INEGI General Directorates



Questions?

Juan.Munoz@inegi.org.mx
Dr. Juan Muñoz López



INSTITUTO NACIONAL
DE ESTADÍSTICA Y GEOGRAFÍA



Conociendo México

01 800 111 46 34

www.inegi.org.mx

atencion.usuarios@inegi.org.mx



@inegi_informa



INEGI Informa



**INSTITUTO NACIONAL
DE ESTADÍSTICA Y GEOGRAFÍA**

