

«Big data» versus «small data»: the case of “gripe” (flu) in Spanish

Antonio Moreno-Sandoval (Department of Linguistics UAM & IIC)

Esteban Moro (Department of Mathematics UC3M & IIC)

[The starting point]

“ ‘Big data hubris’ is the often implicit assumption that big data are a substitute for, rather than a supplement to, traditional data collection and analysis.”

Lazer, Kennedy, King and Vespignani: The Parable of Google Flu: Traps in Big Data Analysis.
Science, vol. 343, 14 March 2014

[The Google Flu Tool (GFT) failure]

→ *Nature* reported in Feb 2013 that GFT was predicting more than double of doctor visits for influenza than the USA medical authorities.

- Method: **counting flu-related searches**
- GFT 2009:
 - “big data were overfitting the small number of cases”
 - “GFT was part flu detector, part winter detector”
- GFT 2013, again, “has been persistently overestimating flu prevalence for a much longer time”

[The error cause, according to Google]

1. **Big data overestimation**
2. **Algorithm dynamics**, which pollute and manipulate data by expanding rumors and trending topics.

[What is Big Data?]

- The Four V's of Big Data by IBM:
 1. **Volume**: scale of data (2.3 trillion Gb created each day)
 2. **Velocity**: analysis of streaming data
 3. **Variety**: different forms of data (text, video, healthcare records, etc.)
 4. **Veracity**: uncertainty of data

[Wikipedia definition]

Big data is a broad term for **data sets** so **large and complex** that traditional data processing applications are inadequate.

The term often refers simply to the use of **predictive analytics** to extract value from the **(unstructured) data**.

[What is 'Big LANGUAGE Data']

- 30 billion pieces of content are shared on Facebook every month
- 4 billion hours of video on YouTube each month
- 400 million tweets are sent per day by about 200 million monthly active users.
- In dozens of different languages.

Can we linguistically process big language data?

- Our experiment: to replicate Google's using Twitter messages in Spanish that were:
 1. geolocalized = Spain
 2. included the word 'gripe'
 3. time span: from Jan 2012 to Aug 2014→ we collected our Spanish Flu Corpus on Twitter

[Spanish Flu Corpus in Twitter]

- 2759 tweets (including RT)
- 327072 'words'
- 140 weeks
- sent from locations in Spain

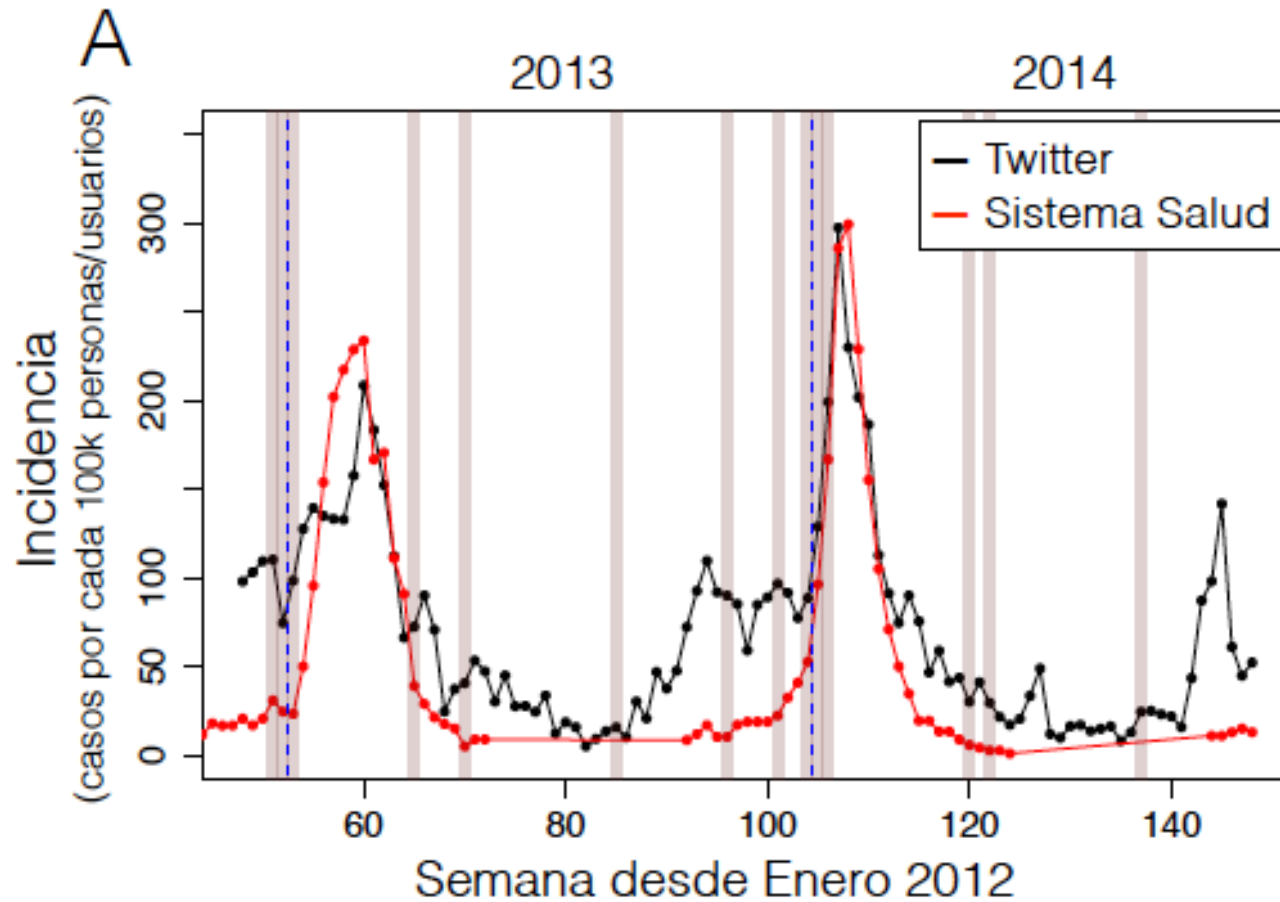
[Procedure]

- **Hypothesis**: the increased number of messages with **GRIPE** is a predictor of an approaching peak of cases
- **Verification**: to check the prediction against the reported cases by the Spanish Health System (real data sent by doctor in health centers and hospitals)

[Elimination of noise in data]

- noise: institutional or press messages:
 - “100.000 personas aún no han podido vacunarse contra la gripe”
- we removed all messages with a URL
- ‘good data’:
 - “estoy en la cama con gripe”

[Results]



[Discussion]

- The messages on Twitter in Spanish also magnify the real cases of flu, as Google (<http://www.google.org/flutrends/es/#ES>.)
- GFP is based on **flu-related searches**.
They compared the query counts with the flu surveillance systems

[The problem: VERACITY in data]

- Analysis of the Flu Corpus to discover what factors contribute negatively to the prediction.
- Distinguishing **Good Data** from **Bad Data**

[Bad data]

■ Figurative and humorous language:

a joke:

“No sé si es un constipado **virulento** o una **GRIPE virurápida**.”

‘ I don’t know if this is a virulent cold or a virufast flu’

- separating real cases from figurative ones requires a good management of pragmatic aspects (intention, irony, metaphor...)

[More Bad Data]

- when the text contains ‘gripe’ but the speaker doesn’t suffers from it:

“Acabo de ver un anuncio de Gelocatil gripe y me he acordado de @...”

‘I just saw an ad of XXX flu and I have remembered @’

→ all those wrong cases contribute to overestimation.

[Some refinements]

- Finding synonyms and variants of ‘gripe’ in the Twitter corpus: ‘gripazo’ ‘griposo’
- Detecting sequences that indicate that someone has contracted this disease:
 - “con la gripe en casa”
 - “menudo gripazo he pillado”
 - “con tos y moqueando”
 - “la gripe me mata”

→ THERE ARE VERY FEW LEXICAL AND SYNTACTIC PATTERNS IN THE TWITTER CORPUS

collocations 'tener' & 'gripe'

- Si **tienes** <gripe > lo mejor es descansar.
- además **tengo** un < gripazo > de cuidado
- Seguramente **tenga** <gripe >
- Odio **tener** <gripe > en temporada de calor
- **tengo** un < gripazo > flipante
- que **tengo** el < gripazo > padre
- Creo que **tengo** <gripe > post-estrés

[Now, the Small Data]

- Analysing 'gripe' in a medical corpus (MultiMedica) with The Sketch Engine: just **341** occurrences against **2759**

'gripe' collocates with:

'virus', 'brote', 'caso' 'estación', 'azote',
'epidemia', 'vacuna' 'temporada'

'padecer' 'sufrir' 'tratar' 'cambiar' 'superar'

['gripe' Word Sketch

gripe (noun)

MULTIMEDICA en Español freq = 341 (74.4 per million)

<u>object of</u> <u>5</u> 0.2	<u>n_modifier</u> <u>44</u> 0.5	<u>modifies</u> <u>201</u> 2.2	<u>v_o</u> <u>11</u> 0.7
sospechar+se <u>1</u> 8.21	2007 <u>3</u> 10.58	temporada <u>16</u> <u>17</u> 11.13	gbs <u>2</u> 10.71
padecer <u>1</u> 4.99	20082009 <u>2</u> 10.48	estación <u>1</u>	escara <u>1</u> 9.14
tratar <u>1</u> 4.92	1918 <u>2</u> 10.3	virus <u>94</u> <u>101</u> 10.33	resfriado <u>1</u> 8.89
sufrir <u>1</u> 4.27	porcino <u>2</u> 9.91	antígeno <u>3</u> cepa <u>4</u>	16 <u>1</u> 8.33
tener <u>1</u> 1.15	antigénicamente nuevo <u>1</u> 9.51	pandemia <u>8</u> 9.85	sífilis <u>1</u> 6.44
	1973 <u>1</u> 9.51	brote <u>9</u> <u>13</u> 8.69	vacuna <u>1</u> <u>2</u> 4.64
	20072008 <u>1</u> 9.51	epidemia <u>4</u>	virus <u>1</u>
<u>subject of</u> <u>50</u> 2.6	2008 <u>1</u> 9.39	epidemiología <u>3</u> 8.45	fiebre <u>1</u> <u>3</u> 4.64
ingresar <u>1</u> 8.85	pandémico <u>1</u> 9.25	vacuna <u>14</u> 8.26	infección <u>1</u>
ee <u>1</u> 8.13	verdad <u>1</u> 9.14	azote <u>1</u> 7.33	mortalidad <u>1</u>
influir <u>1</u> <u>2</u> 6.77			

[Conclusions]

- Our study supports Lazer et al. 2014: “instead of focusing on a ‘big data revolution,’ perhaps it is time we were focused on an **‘all data revolution,’** where we recognize that the critical change in the world has been innovative analytics, **using data from all traditional and new sources,** and providing a deeper, clearer understanding of our world.”

[Conclusions (2)]

- In terms of corpora, small but well selected collections of linguistic data should be combined with large repositories from internet and social networks, since sometimes **'small data' offer information that is not inferred from 'big data'**.