

Text Mining II - Máster Big Data Analytics (2015 - 2016)

Exploración del dataset HispaTweets

Dr. Paolo Rosso

Casi-casi-Dr. Francisco Rangel ;-)



UNIVERSITAT
POLITÈCNICA
DE VALÈNCIA

Objetivos de la exploración de datos

- Entender propiedades de los datos.
- Encontrar patrones.
- Depurar.
- Comunicar resultados.

Descarga del dataset HispaTweets (Fabra et al., 2016)

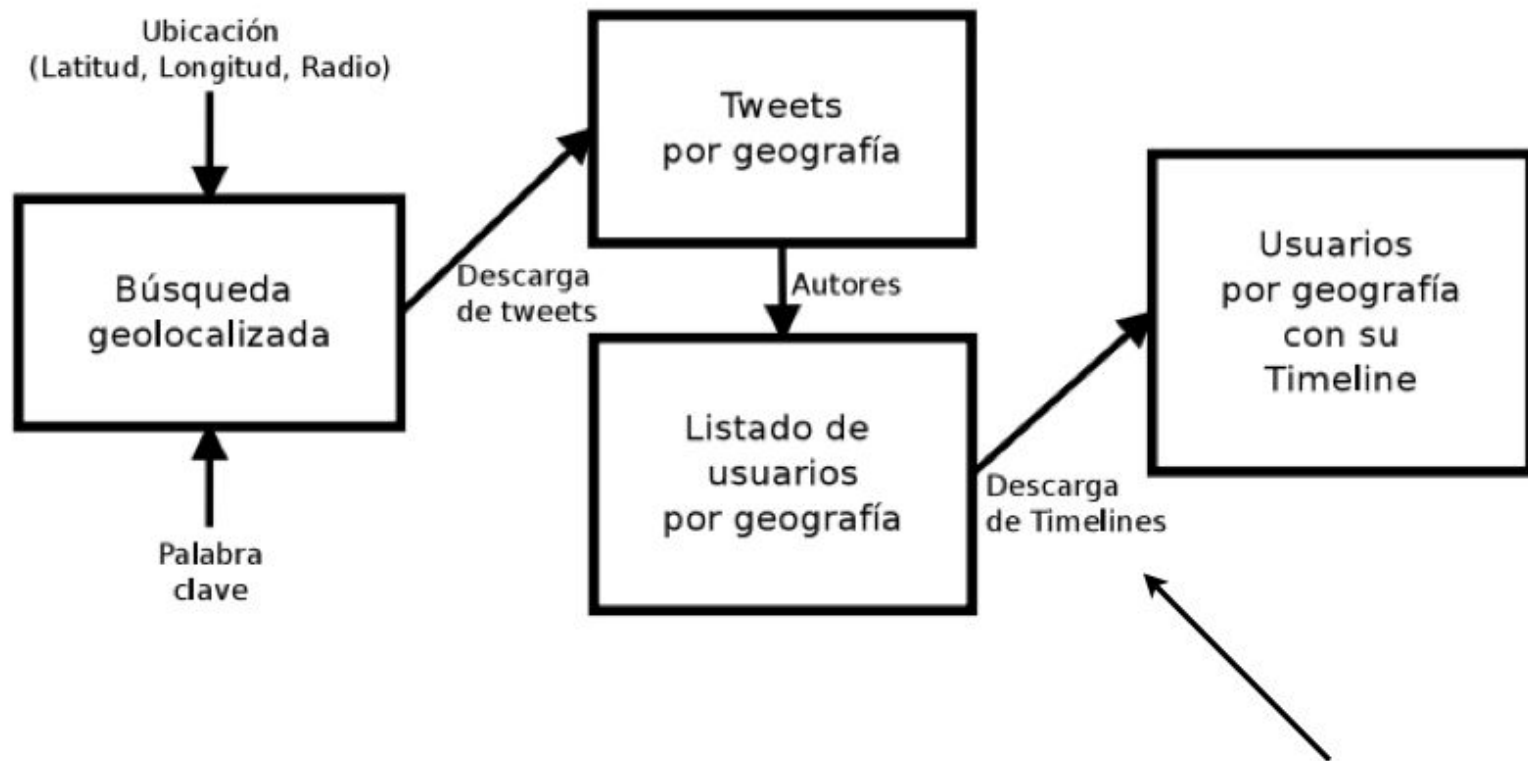
<https://s3.amazonaws.com/cosmos.datasets/hispatweets.zip>

Características del dataset

- Obtenido de Twitter.
- Colección de decenas/cientos de miles de autores.
- Cientos de tuits por autor.
- Gran variedad de temas.
- ¿Gran tamaño? 700Mb comprimido, 13Gb descomprimido.
- Dificultad para etiquetar la información:
 - Variedad del lenguaje por posición geográfica.
 - Sexo por nombre (mujeres, hombres y ¿perfiles institucionales?)
 - Personas reales vs. robots (chatbots)

*¿Big
Data?*

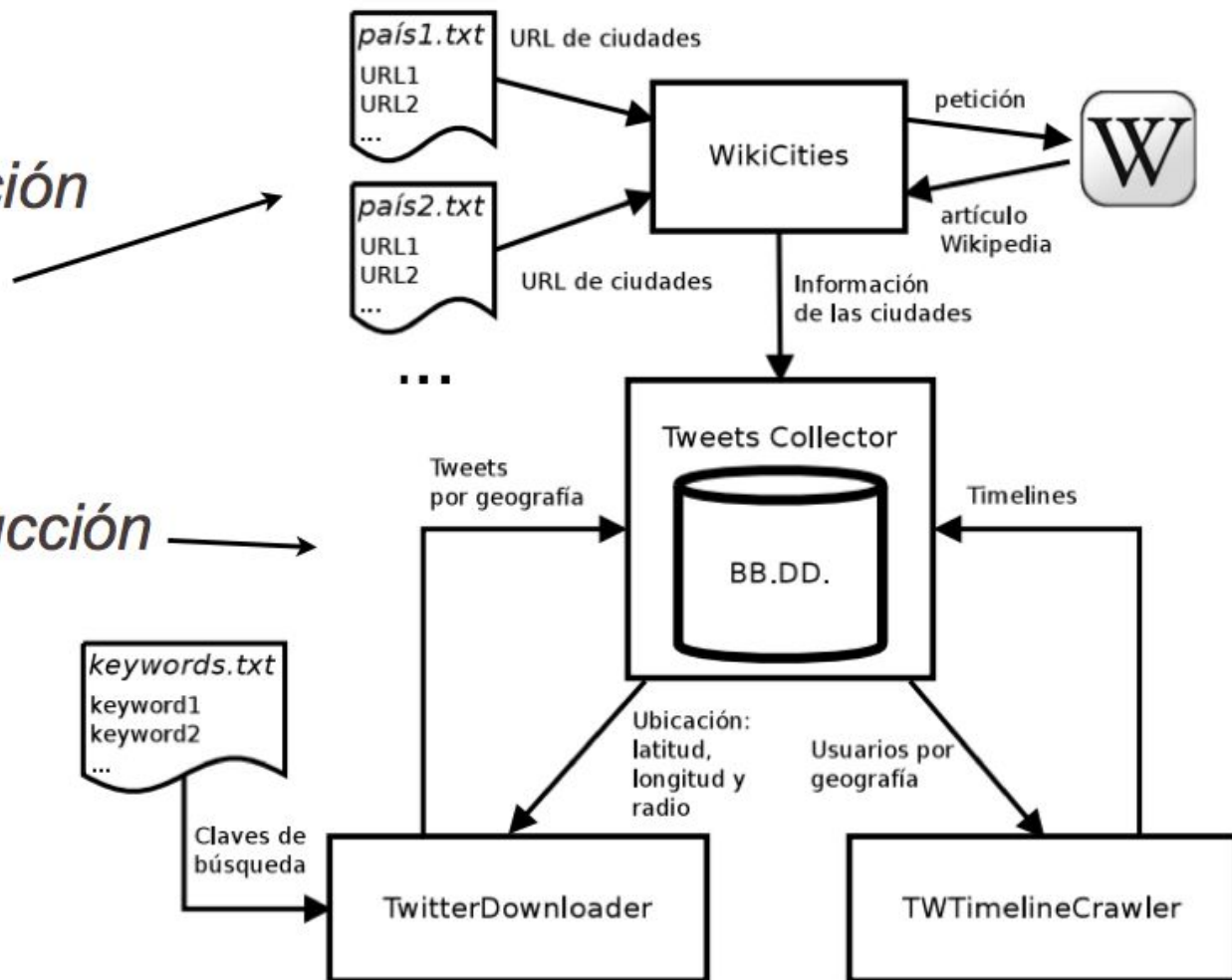
Proceso de construcción



hasta los últimos 1000 tweets de sus respectivas *timelines*.

*Proceso de obtención
de coordenadas
geográficas*

*Proceso de construcción
del corpus*

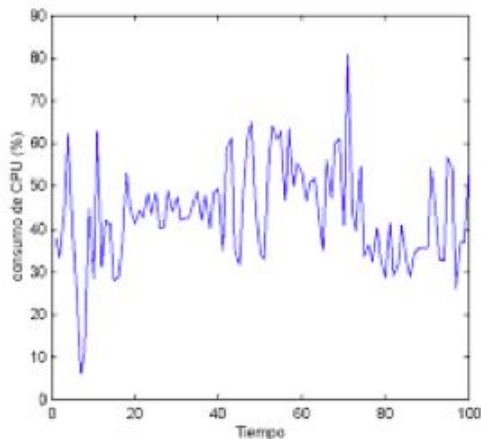


Criterio de selección

Depuración, refinamiento y filtrado del corpus:



Geográfico



Temporal



Usuarios

Etiquetado del sexo

Basado en diccionario de nombres propios:

Mujeres	24.429
Hombres	25.949
<i>Total</i>	<i>50.378</i>



Formato de los ficheros

Una vez descomprimido el dataset, la estructura es la siguiente:

- Un par de ficheros de verdad: training.txt y test.txt. El formato es:
 - id:::variedad:::sexo
- Una carpeta por variedad, un fichero .json por autor:
 - Cada línea del fichero un tuit en formato json.
 - A la derecha un ejemplo ->
 - Se recomienda el uso de jsonlint.com

```
1 {"lang": "es", "quoted_status_id_str": "682341536172871681", "in_reply_
2 {"lang": "und", "quoted_status_id_str": "540503573944336384", "in_reply
3 {"lang": "es", "entities": {"symbols": [], "user_mentions": [{"indices"
4 {"lang": "es", "entities": {"symbols": [], "user_mentions": [{"indices"
5 {"lang": "es", "entities": {"symbols": [], "user_mentions": [], "hashta
6 {"lang": "es", "entities": {"symbols": [], "user_mentions": [], "hashta
7 {"lang": "es", "entities": {"symbols": [], "user_mentions": [{"indices"
8 {"lang": "es", "entities": {"symbols": [], "user_mentions": [], "hashta
9 {"lang": "es", "entities": {"symbols": [], "user_mentions": [{"indices"
10 {"lang": "en", "entities": {"symbols": [], "user_mentions": [], "hashta
11 {"lang": "es", "entities": {"symbols": [], "user_mentions": [], "hashta
12 {"lang": "pt", "entities": {"symbols": [], "media": [{"indices": [24, 4
13 {"lang": "pt", "entities": {"symbols": [], "user_mentions": [], "hashta
14 {"lang": "pt", "entities": {"symbols": [], "user_mentions": [], "hashta
15 {"lang": "es", "entities": {"symbols": [], "user_mentions": [], "hashta
16 {"lang": "es", "entities": {"symbols": [], "user_mentions": [{"indices"
```

¿Qué nos interesa explorar?

- Número de autores por clase (sexo y variedad del lenguaje)
- Número de tuits por autor.
- Número de tuits por clase.
- Número de palabras por documento / autor / clase.
- Distribución de palabras/documentos/autores por documento/autor/clase...
- Longitud media de tuits, palabras, documentos...por clase.
- Distribución temporal de los tuits, tuit más antiguo, más nuevo, media, desviación...
- Palabras extrañas, frecuentes, comunes...
- ... cualquier información que nos describa el dataset y aporte conocimiento nuevo y valuable.