

Analizando los cables diplomáticos de los EEUU con python

Pablo Angulo Ardoy Juan Duque J. Ángel González Prieto
(David Gómez-Ullate)

Universidad Autónoma de Madrid

Universidad Complutense de Madrid

Vigo, junio de 2012

Wikileaks



- Ofrece **canales seguros** para que aquellas personas con acceso a material comprometido puedan filtrarlo de forma **anónima**.
- Una vez recibido, comprueba la autenticidad del documento.
- Una vez verificado, pone el material en circulación: wikileaks.org, bittorrent ...
- En ocasiones, lo pasa a la prensa para que obtenga mayor difusión o por otros motivos.

Hitos de Wikileaks



- **Diciembre 2006** Wikileaks empieza a operar.
- **Abril 2010** Wikileaks libera *Collateral Murder*.
- **Mayo 2010** Wikileaks ya no es una wiki (editable por todo el mundo).
- **Julio 2010** Wikileaks libera los *Diarios de la Guerra de Afganistán*.
- **Octubre 2010** Wikileaks libera los *Diarios de la Guerra de Irak*.

Los cables diplomáticos filtrados por Wikileaks

- **Julio 2010** Wikileaks pone en circulación un **fichero encriptado** que contiene **cables diplomáticos** a modo de "*seguro de vida*".
- **Noviembre 2010** *El País*, *Le Monde*, *Der Spiegel*, *The Guardian*, y *The New York Times* comienzan a publicar cables, muy poco a poco y editados para no poner en peligro a las personas mencionadas en los cables.
- **Febrero 2011** Periodistas de *The Guardian* publican la clave de encriptación del fichero, pensando que ya no era útil.
- **Agosto 2011** *Der Freitag* observa que tanto el fichero como la clave están online, con lo que es posible obtener todos los cables sin editar.

El conjunto de cables

- 251,287 cables.
- 87 MB de info en las cabeceras, 1.5GB en el cuerpo.
- Desde 1966 a 2010. (1000 cables hasta el año 2000, 56813 sólo en 2009)
- Al principio muy poco tráfico, poco a poco se van uniendo más embajadas a la red.
- Cada cable contiene una **cabecera**, con información sobre remitente, destinatarios... y un **cuerpo**, con título, resumen, etiquetas y cuerpo del mensaje.

Cabecera de un cable

VZCZCXRO0967

RR RUEHIK RUEHLN RUEHPOD RUEHSK

RUEHVK RUEHYG

DE RUEHMO #0252/01 0331531

ZNR UUUUU ZZH

R 021531Z FEB 09

FM AMEMBASSY MOSCOW

TO RUEHC/SECSTATE WASHDC 1740

INFO RUCNCIS/CIS COLLECTIVE

RUEHXD/MOSCOW POLITICAL COLLECTIVE

RUEHZG/NATO EU COLLECTIVE

RUEAIIA/CIA WASHDC

RUEKJCS/SECDEF WASHDC

RUEKDIA/DIA WASHDC

RUEKJCS/JOINT STAFF WASHDC

RHEHNSC/NSC WASHDC

Cuerpo del mensaje

UNCLAS SECTION 01 OF 02 MOSCOW 000252

SENSITIVE

SIPDIS

E.O. 12958: N/A

TAGS: PGOV PREL PINS MARR MASS RS AF

SUBJECT: GOR SUPPORTS ISAF TRANSIT OF GOODS TO AFGHANISTAN,
SEEKS EXPERTS MEETING

REF: A. STATE 6471

¶B. 08 MOSCOW 3655

¶C. SZPILA-EMB MOSCOW E-MAIL JANUARY 29

¶1. (SBU) Summary and Action Request: In response to ref A
dipnote, the MFA has indicated that, provided all usual

...

of Russia. End Summary and action request.

Transit Deal Is Complete, Now Let's Talk About It

...

¿Qué significa todo esto?

- **cablsearch**, **look4leaks** y otros sitios web nos dicen qué significan las etiquetas.
- **Wikipedia**: información general sobre la diplomacia.
- Observar cables, conjeturar lo que significa cada cosa, ignorar lo que no necesitas.
- Muchos acrónimos y jerga se entienden si lees unos cuantos cables de la misma embajada.

Cabecera explicada

VZCZCXRO0967 ???

RR: Prioridad del cable

RR RUEHIK RUEHLN RUEHPOD RUEHSK **RUEH??:** Para diversas embajadas

RUEHVK RUEHYG

RUEHMO: Embajada de Moscú (remitente)

DE RUEHMO #0252/01 0331531

#0252/01: Identificador del cable

ZNR UUUUU ZZH ???

R 021531Z FEB 09 02 de Febrero de 2009, 15:31 GMT

FM AMEMBASSY MOSCOW

FM: From (remitente)

TO RUEHC/SECSTATE WASHDC 1740

TO: Destinatario

INFO RUCNCIS/CIS COLLECTIVE

RUEHC/SECSTATE: Secretaría de Estado

RUEHXD/MOSCOW POLITICAL COLLECTIVE

RUEHZG/NATO EU COLLECTIVE

RUEAIIA/CIA WASHDC

INFO: Otros destinatarios

RUEKJCS/SECDEF WASHDC

RUEKDIA/DIA WASHDC

RUEKJCS/JOINT STAFF WASHDC

RHEHNSC/NSC WASHDC

Anatomía de un cable

Cuerpo explicado

UNCLAS SECTION 01 OF 02 MOSCOW 000252

UNCLAS: Unclassified

SENSITIVE

SENSITIVE: Pero un tema sensible

SIPDIS

SIPDIS: Proveniente de SIPRNet

E.O. 12958: N/A ???

(A pesar de que no es secreto)

TAGS: PGOV PREL PINS MARR MASS RS AF

TAGS: Sistema de clasificación de cables

(Hablaemos sobre ello más adelante)

SUBJECT: GOR SUPPORTS ISAF TRANSIT OF GOODS TO AFGHANISTAN,

SUBJECT: Tema del cable

SEEKS EXPERTS MEETING

REF: A. STATE 6471

REF: El cable menciona otros cables, emails, etc

¶B. 08 MOSCOW 3655

¶C. SZPILA-EMB MOSCOW E-MAIL JANUARY 29

¶1. (SBU) Summary and Action Request: In response to ref A

dipnote, the MFA has indicated that, provided all usual

Summary: Resumen

...

of Russia. End Summary and action request.

Sigue el cuerpo del cable, con secciones delimitadas de formas diversas.

Transit Deal Is Complete, Now Let's Talk About It

Usa acrónimos, algunos explicados y otros no (véase GOR: Government of Russia).

...

TAGS: etiquetas de países

- AA** Aruba
- AC** Antigua and Barbuda
- AE** United Arab Emirates
- AF** Afghanistan
- AG** Algeria
- AJ** Azerbaijan
- AL** Albania
- AM** Armenia
- AN** Andorra
- AO** Angola

...

TAGS: etiquetas de temas

PREL External Political Relations

PGOV Internal Governmental Affairs

PHUM Human Rights

ECON Economic Conditions

PTER Terrorists and Terrorism

PINR Intelligence

KPAO Public Affairs Office

OIIP International Information Programs

...

TAGS: otras etiquetas

Programas

KDEM Democratization

KCRM Criminal Activity

KMDR Media Reaction Reporting

...

Organismos

IAEA International Atomic Energy
Agency

NATO North Atlantic Treaty Organization

...

Problemas antes de empezar a trabajar

- Cables incompletos (49082 cables, a partir de 2007 todos están completos)
- Inconsistencias, palabras mal deletreadas ...

Tags Presentes casi siempre, a veces mal deletreadas, pero en general muy útil.

Subject Presente casi siempre.

Summary Aprox la mitad de los cables.

Ref Aprox la mitad de los cables. Referencias a cables en multitud de formatos, pero también a emails, faxes...

Operaciones con texto

- Comprobar si una cadena está en el texto:

```
1 def has_summary(texto):  
2     return 'summary' in texto
```

- Comprobar si una cadena está en el texto (ignorando mayúsculas y minúsculas):

```
1 def has_summary(texto):  
2     return 'summary' in texto.lower()
```

Expresiones regulares

Definimos un patrón, y buscamos texto que encaja en ese patrón:

- \w** Acepta cualquier carácter alfanumérico.
- \s** Acepta cualquier espacio en blanco.
 - Acepta cualquier carácter.
- []** Agrupar varios caracteres (ej: `x'[\s,]'` espacio en blanco, o coma)
- +** Una o más apariciones (ej: `'\s+'` uno o más espacios en blanco)
- *** Cero o más apariciones
- ?** Cero o una apariciones
- ()** Guarda parte de la cadena para uso posterior.

Expresiones regulares aplicadas a cables

- Contiene la palabra 'end', seguida de uno o más espacios, seguida la palabra 'summary' (ignorando mayúsculas y minúsculas):

```
1 import re
2 END_SUMMARY = re.compile(r'end\s+summary', re.
   IGNORECASE)
3 def has_end_summary(texto):
4     return bool(END_SUMMARY.search( texto))
```

Expresiones regulares aplicadas a cables

- Buscar la línea que contiene las etiquetas:

```
1 >>> TAGS = re.compile(r'TAGS:?\s+(.*)', re.  
    IGNORECASE)  
2 >>> m = TAGS.search(cuerpo)  
3 >>> print m.group(1)  
4 'PGOV PREL PINS MARR MASS RS AF'
```

Expresiones regulares aplicadas a cables

- Separar las etiquetas, que pueden estar separadas por espacios, comas, punto y coma, coma y espacio...:

```
1 >>> sep = re.compile(r'[\s,;]+')
2 >>> sep.split('PGOV PREL PINS MARR MASS RS AF')
3 ['PGOV', 'PREL', 'PINS', 'MARR', 'MASS', 'RS', 'AF']
4 >>> sep.split('PGOV, PREL, PINS, MARR, MASS, RS, AF'
5 )
   ['PGOV', 'PREL', 'PINS', 'MARR', 'MASS', 'RS', 'AF']
```


Extraer información estructurada de la wikipedia (II)

- Para descargar el html de una url:

```
1 import urllib2
2 def download(url):
3     httpFile = urllib2.urlopen(urllib2.Request(url))
4     httpBody=httpFile.read()
5     httpFile.close()
6     return httpBody
```


Extraer información estructurada de la wikipedia (III)

- Extraer la capital de un país de la página de la wikipedia que habla de ese país.

```
1 def getCapital(link):
2     countryPage=download(baseUrl+link)
3     result=re.search('Capital.*\n+.*\n*<td
4         >(?:<[^>]+>)?(?:<[^<]*<',countryPage)
5     if result is None:
6         return None
7     return result.group(1)
```

Extraer información estructurada de la wikipedia (IV)

- Extraer una lista de países de la wikipedia, con el link a la página de la wikipedia sobre el país.

```
1 listOfCountriesURL=baseURL+"wiki/  
   List_of_sovereign_states"  
2 countriesPage = download(listOfCountriesURL)  
3 COUNTRIES = re.compile('<span id=.*></span  
   >.*&#160;.*<a href="(.)" title=".*">(.)</a>')  
4 for link,country in COUNTRIES.findall(countriesPage)  
   :  
5     capital=getCapital(link)  
6     ...
```

Índice

1 Los datos

- Origen de los datos
- Anatomía de un cable
- Etiquetas
- Preprocesado

2 Análisis de comunidades

3 Análisis temático

- Análisis de etiquetas
- Clustering
- SVD

Networkx

NetworkX

[NetworkX Home](#) | [Download](#) | [Developer Zone](#) | [Documentation](#) | [Blog](#) >

Random Geometric Graph

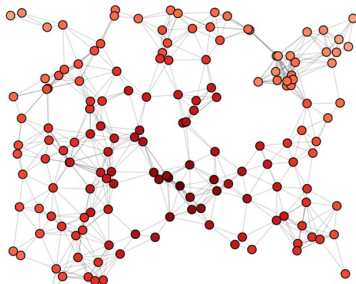


Figure: Paquete networkx (Python)

Networkx

- Estructura de red = diccionario

```
1 import networkx as nx
2
3 #Generacion de la red de atributo name
4 G = nx.Digraph(name = 'wikileaks')
5
6 #G.graph devuelve {'name': 'wikileaks'}
```

Nodos

- Campo remitente
- Atributos
 - longitud
 - latitud
 - cables totales emitidos
 - cables en red

Nodos: lectura de atributos

```
1
2  # Acceso a los atributos de un nodo
3  #Totalidad de los nodos de la red con sus atributos
4  G.nodes(data = True)
5  [(Nodo1, {...}), (Nodo2, {...}),...]
6
7  #Diccionario de atributos de un nodo
8  G[Nodo1]
9  {...}
10
11 #Valor del atributo latitud de un nodo
12 G[Nodo1]['latitud']
13 ...
```

Links

- Campos to e info.
- Atributos
 - lista to
 - lista info
 - peso = len(to U info)
 - fechas de emisión.

Visualizando la red: Gephi

- Formatos: gexf, gml, GraphML,...
- Escritura en networkx: `nx.write_gexf`, `nx.write_gml`,
`nx.write_graphml`,...

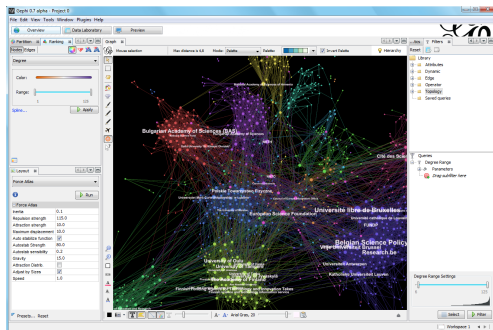


Figure: Gephi como herramienta de visualización

El problema de la modularidad

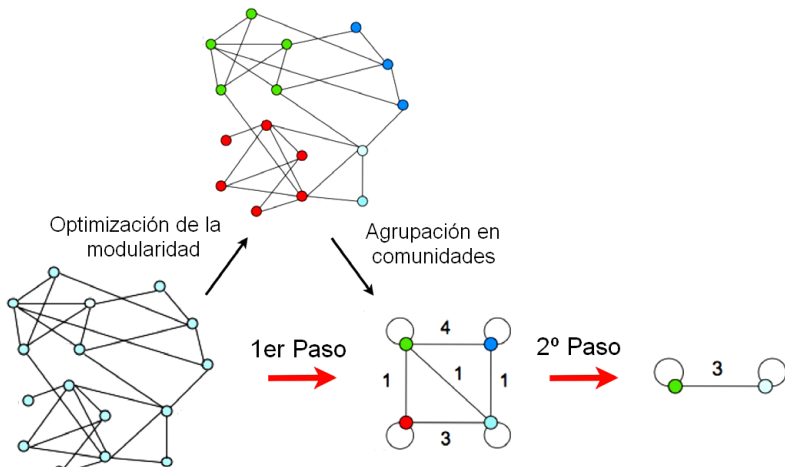
- Diferencia entre el número de links entre las comunidades menos el número que hubiera si estuvieran generados al azar

$$Q = \frac{1}{2m} \sum_{ij}^N \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j)$$

donde:

- k_i es el grado del nodo i .
- m es el número de links
- c_i es la comunidad del nodo i .
- $\delta(c_i, c_j) = 1$ si $c_i = c_j$ y 0 en cualquier otro caso.

Algoritmo para el cálculo de la modularidad: Método de Blondel-Guillaume-Lambiotte-Lefebvre



Ejemplo de ejecución

- `http://perso.crans.org/aynaud/communities/:community.py`
- `import community`
- `import networkx as nx`
- `G=nx.read_gml("karate.gml")`
- `partition = community.best_partition(G)`
- `modularity = community.modularity(partition, G)`

Ejemplo de ejecución

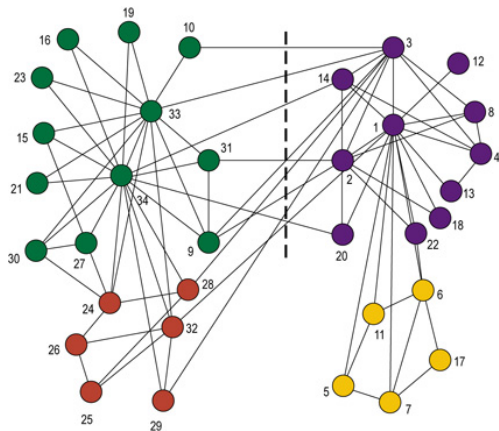


Figure: Ejemplo de detección de comunidades: Red de Zachary del club de Karate.

Resultados

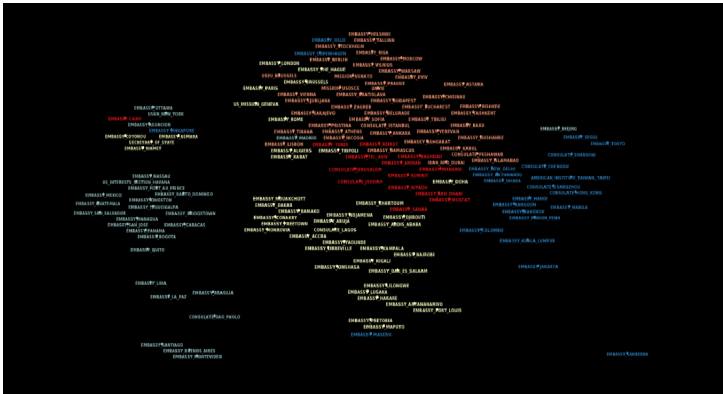


Figure: Modularidad y estructura de comunidades de la red de wikileaks: año 2009

Índice

1 Los datos

- Origen de los datos
- Anatomía de un cable
- Etiquetas
- Preprocesado

2 Análisis de comunidades

3 Análisis temático

- Análisis de etiquetas
- Clustering
- SVD

Fundamentos

Idea

Realizar un análisis del corpus basado en la temática de los mensajes.

Problema: Inferir el tema de un cable.

Solución: Utilizar las tags como marcadores semánticos.

Ventajas

- **Objetividad:** Las tags son fijadas por el creador del mensaje.
- **Especialización:** Las tags han sido creadas por expertos diplomáticos.

